

PCT

WORLD INTELLECTUAL PROPERTY ORGANIZATION  
International Bureau



39

INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

<b>(51) International Patent Classification <sup>6</sup> :</b> <b>C07K 1/00, 14/47, G06F 17/50</b>	<b>A2</b>	<b>(11) International Publication Number:</b> <b>WO 99/42474</b> <b>(43) International Publication Date:</b> 26 August 1999 (26.08.99)
<b>(21) International Application Number:</b> PCT/US99/03692 <b>(22) International Filing Date:</b> 19 February 1999 (19.02.99)  <b>(30) Priority Data:</b> 60/075,466 20 February 1998 (20.02.98) US  <b>(71) Applicant:</b> GENOME DYNAMICS, INC. [US/US]; 18247D Flower Hill Way, Gaithersburg, MD 20879-0921 (US).  <b>(72) Inventors:</b> MICHAELS, George, S.; 13516 Sloan Street, Rockville, MD 20853 (US). MIKELSAAR, Raik-Hiio; Jakobsoni Street 11-4, EE2400 Tartu (EE). FELDMANN, Richard, J.; 17800 Mill Creek Drive, Derwood, MD 20855-1019 (US).  <b>(74) Agents:</b> DRIVAS, Dimitrios, T. et al.; White & Case LLP, 1155 Avenue of the Americas, New York, NY 10036 (US).		<b>(81) Designated States:</b> AU, BR, BY, CA, CN, CZ, HU, IL, IN, JP, KR, MX, NZ, PL, RU, SK, European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE).  <b>Published</b> <i>Without international search report and to be republished upon receipt of that report.</i>
<b>(54) Title:</b> METHOD FOR DESIGNING DNA-BINDING PROTEINS OF THE ZINC-FINGER CLASS  <b>(57) Abstract</b>  The invention is directed to the design of DNA-binding proteins (DBP's) with capabilities of binding to any predetermined target double-stranded DNA sequence. Disclosed are the rules for design of the proteins; an algorithm for screening for the optimal DBP's; a computer system employing the rules and the algorithm; general formulae encompassing the proteins; and methods of use of the proteins.		

**FOR THE PURPOSES OF INFORMATION ONLY**

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AL	Albania	ES	Spain	LS	Lesotho	SI	Slovenia
AM	Armenia	FI	Finland	LT	Lithuania	SK	Slovakia
AT	Austria	FR	France	LU	Luxembourg	SN	Senegal
AU	Australia	GA	Gabon	LV	Latvia	SZ	Swaziland
AZ	Azerbaijan	GB	United Kingdom	MC	Monaco	TD	Chad
BA	Bosnia and Herzegovina	GE	Georgia	MD	Republic of Moldova	TG	Togo
BB	Barbados	GH	Ghana	MG	Madagascar	TJ	Tajikistan
BE	Belgium	GN	Guinea	MK	The former Yugoslav Republic of Macedonia	TM	Turkmenistan
BF	Burkina Faso	GR	Greece	ML	Mali	TR	Turkey
BG	Bulgaria	HU	Hungary	MN	Mongolia	TT	Trinidad and Tobago
BJ	Benin	IE	Ireland	MR	Mauritania	UA	Ukraine
BR	Brazil	IL	Israel	MW	Malawi	UG	Uganda
BY	Belarus	IS	Iceland	MX	Mexico	US	United States of America
CA	Canada	IT	Italy	NE	Niger	UZ	Uzbekistan
CF	Central African Republic	JP	Japan	NL	Netherlands	VN	Viet Nam
CG	Congo	KE	Kenya	NO	Norway	YU	Yugoslavia
CH	Switzerland	KG	Kyrgyzstan	NZ	New Zealand	ZW	Zimbabwe
CI	Côte d'Ivoire	KP	Democratic People's Republic of Korea	PL	Poland		
CM	Cameroon	KR	Republic of Korea	PT	Portugal		
CN	China	KZ	Kazakhstan	RO	Romania		
CU	Cuba	LC	Saint Lucia	RU	Russian Federation		
CZ	Czech Republic	LI	Liechtenstein	SD	Sudan		
DE	Germany	LK	Sri Lanka	SE	Sweden		
DK	Denmark	LR	Liberia	SG	Singapore		
EE	Estonia						

## METHOD FOR DESIGNING DNA-BINDING PROTEINS OF THE ZINC-FINGER CLASS

### BACKGROUND OF THE INVENTION

A superfamily of eukaryotic genes encoding potential nucleic-acid-binding proteins contains zinc-finger (ZF) domains of the Cys<sub>2</sub>-His<sub>2</sub> (C<sub>2</sub>H<sub>2</sub>) class. Proteins that have these characteristic structural features play a key role in the regulation of gene expression [1-4]. Sequence comparisons, mutational analyses, and a recent crystallographic investigation have revealed that each finger domain, as a rule, interacts with the major groove of B-form DNA through contacts with some or all three base pairs within a DNA triplet. These base-specific interactions are mediated through amino acid (AA) side chains at specific positions in the  $\alpha$ -helical region [5-10] of the protein domain.

Although the AA sequences of more than 1,300 ZF motifs have been identified, the exact DNA-binding sites are known only for a few proteins. The available information on DNA contact regions concerns mainly guanine-cytosine-rich strands [5-9] and fewer adenine-thymine-rich sites [11,12]. On the basis of experimental data, the first proposals for rules relating ZF sequences to preferred DNA-binding sites have been made [13,14]. However, no general rules for ZF protein-DNA recognition have been proposed. This is likely due to the fact that neither computer modeling [2,3,5] nor crystallographic analysis [7] have provided enough information on the overall structural variety in the ZF-DNA contact region.

Using physical atomic-molecular models to characterize the steric conditions in the specific contact positions for different ZF-DNA interactions, an objective of the work leading to the present invention was to determine a set of general rules for ZF-DNA recognition for the C<sub>2</sub>H<sub>2</sub> class of ZF domains. Once this objective had been reached, the work of the invention plan was to develop an algorithm, and a computer system using the algorithm, to design effective zinc-finger DNA-binding polypeptides. The achievement of these goals represents a major advance of knowledge in the field, knowledge characterized by the

disclosures of Rebar, et. al. and Beerli, et. al. [15,16]. These two disclosures are concerned with the selection, using the phage display system, of specific zinc fingers with new DNA-binding specificities. On the other hand, the present disclosure is concerned with the design of DNA-binding proteins for any given DNA sequence.

## SUMMARY OF THE INVENTION

The invention is directed to the design and specification of DNA-binding proteins binding via  $C_2H_2$  zinc-finger motifs (DBP's or, individually, a DBP). On the basis of the studies described herein, general rules for optimizing such binding have been determined, and a formula describing the class of DBP's having optimal DNA-binding properties has been constructed. Furthermore, a program has been developed, based on the rules, which affords the design of DBP's with such high binding affinity for any given DNA sequence. Lastly, rules have been determined for the design of DBP's which, while not having optimal binding, do have significant and useful DNA-binding properties.

## BRIEF DESCRIPTION OF THE DRAWINGS

Figure 1 depicts the alignment of ZF domains in various known DBP's.

Figure 2 is a schematic representation of the interaction between a target DNA triplet and a single ZF domain.

Figure 3 is a schematic representation of the interaction between a target DNA string of 9 bases and a three-domain DBP.

Figure 4 is a block flow diagram of the computer system by which the instant DBP design process is implemented.

Figure 5 is a block flow diagram wherein the Computer Program block (2) of Figure 4 is further broken down.

Figure 6 is a block flow diagram wherein the Process Genome into Blocking Fragment Files block (2) of Figure 5 is further broken down.

Figure 7 is a block flow diagram wherein the Design DBP's for a Genome block (3) of Figure 5 is further broken down.

Figure 8 is a block flow diagram wherein block (22) of Figure 7 is further broken down.

Figure 9 is a block flow diagram wherein block (24) of Figure 7 is further broken down.

Figure 10 shows the distribution of binding strengths of acceptable 9-finger DBP's across the yeast genes analyzed.

Figure 11 shows the values of the binding energies of the acceptable 9-finger DBP's found for the yeast genes analyzed.

Figure 12 shows the distribution of DBP subsite (spurious) binding energies across the yeast genes analyzed.

Figure 13 shows, in nonlogarithmic fashion, the distribution depicted in Figure 12.

Figure 14 shows the ratios of binding energy to subsite (spurious) binding energy, across the yeast genes analyzed, for the acceptable 9-finger DBP's.

Figure 15 shows the values of the spurious binding energies for each of the 27-base-pair (bp) frames of the 300-bp promoter region of yeast gene YAR073.

Figure 16 shows the ratios of binding energy to subsite (spurious) binding energy for each of the 27-base-pair (bp) frames of the 300-bp promoter region of yeast gene YAR073.

Figure 17 shows the distribution of sizes of acceptable DBP's across the *C. elegans* genes analyzed.

Figure 18 shows the ratios of binding energy to subsite (spurious) binding energy, across the *C. elegans* genes analyzed, for the acceptable DBP's.

## DETAILED DESCRIPTION OF THE INVENTION

The general rules governing the binding of  $C_2H_2$  ZF motifs to DNA were developed by using a combination of the database analysis of the homologies between 1,851 possible ZF domains and physical molecular modeling of the interaction of a DBP model with a DNA model containing all 64 possible base-pair triplets. The DBP model approximates the size and shape of a half-gallon jug of milk. The DNA model is approximately four feet long and one foot in diameter. The axis of the DNA model is horizontal and can be rotated to observe each of the 64 base-pair triplets. By moving the DBP model in and out with respect to the DNA model one can observe the amino acid and nucleic acid contacts.

Although the following description details the scientific precedents of this invention, the completeness of the rule set governing the DBP-DNA interaction could have only been obtained by the continual, derivative interplay of data base analysis and physical modeling during the invention period. Observations as to the conservation and variability of amino acids at various places in the ZF motif were embodied, first, by constructing a physical model of the ZF motif and, then, by physically modeling the interaction of a specific DBP with a designated DNA bp triplet. The physical modeling indicated patterns of amino acid and nucleic acid interaction which led to further analysis of the database. Iterations of this interplay between database analysis and physical modeling enabled conceptual refinement and expansion of the nature of contact patterns. As these patterns emerged, systematic variation of the amino acids in the ZF motif was undertaken for each of the 64 base-pair triplets. The physical modeling of the interaction between a DBP and DNA was efficient because alternative amino acids could be easily introduced into the ZF motif and the resulting protein physically modeled against the DNA. Hydrogen bonding, and water and hydrophobic contacts could then be modeled, clearly determined and counted very quickly. From this physical modeling a general set of rules was developed which incorporates criteria for the design of DBP's that specifically interact with DNA.

The utility of ZF sequence analysis and alignment is illustrated by Figure 1. The TFIIIA protein is widely used as a model for ZF proteins both in terms of physical measurement and modification and theoretical data analysis. For each of the nine zinc-finger domains the TFIIIA amino acid sequence in this figure has been aligned so that the zinc-binding amino acids, the two cysteines (CYS) and the two histidines (HIS), are aligned in four columns. In order to achieve this alignment dashes must be inserted into the sequence at

various places to provide for domains which have additional amino acids. The same type of alignment has been done for ZF protein MKR2 and the Kruppel proteins. The MKR2 sequence alignment is very compact; there is no need for any insertions, since all of its ZF domains are of the same size. Compared to TFIIIA, MKR2 acts as a much more uniform model for studying the interaction of the amino acids of the protein with the bases of the specific double-stranded DNA. To arrive at the present invention, MKR2 has been used exclusively as the sequence basis for deducing the general rules which govern DBP-DNA interactions.

The crystallographic analysis of a complex containing three ZFs from ZF protein Zif268 and a consensus DNA-binding site helped identify the localization of ZF-B-DNA recognition sub-sites [7]. Because the mutagenesis and sequence investigation results are in accordance with crystal structure data, it is reasonable to expect that the same contact regions also participate in the interaction of other ZF-DNA complexes [5,6,8-10]. Thus, it has been assumed that the following ZF components of a ZF protein play a key role in the anti-parallel DNA reading process: 1) the AA immediately preceding the  $\alpha$ -helical region of the protein; 2) the third residue within the  $\alpha$ -helical region, i.e., that immediately preceding constant leucine; and, 3) the sixth residue of this region, i.e., that immediately preceding invariant histidine.

These components are indicated below as  $Z_3$ ,  $Z_2$  and  $Z_1$ , respectively, in the generalized ZF sequence ( $\alpha$ -helical and  $\beta$ -structural regions are underlined) given in Formula I:

Y/F X C X<sub>2-4</sub> C G/D K/R X F X Z<sub>3</sub> X X Z<sub>2</sub> L X Z<sub>1</sub> H X<sub>3-5</sub> H T/S G/E X<sub>0-2</sub> E K/R P

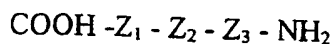
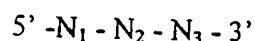
$\beta$ -structural region

$\alpha$ -helix

#### Formula I

wherein X is any amino acid;  $X_{2-4}$  is a peptide 2 to 4 amino acids in length;  $X_{3-5}$  is a peptide 3 to 5 amino acids in length;  $X_{0-2}$  is a peptide 0 to 2 amino acids in length and C, D, E, F, G, H, K, L, P, R, S, T and Y designate specific amino acids according to the standard single-letter code. Pairs of letters separated by "/" indicate that the position can be filled by either of the two specific amino acids designated.

Keeping in mind the above formula, one can envision the formation of antiparallel, trinucleotide-peptide complexes with three (first, second and third) contact positions as follows:



The crystallographic investigation of the Zif268-DNA complex also gave indications of the way the contact groups interact. Pavletich and Pabo [7] concluded that Zif268 forms 11 critical hydrogen bonds (H-bonds) with the bases of the coding DNA strand in the major groove. Two arginine residues in the first contact positions (see the designations of positions above) make H-bonds with the N7 and O6 atoms of the guanine. Three arginine residues hydrogen bond in the same way with guanine in the third contact position. In addition, each arginine residue in this position forms lateral H-bond, salt bridge interactions with carboxylate groups of aspartic acid occurring as the second residue in the  $\alpha$ -helix. The N $\delta$  atom of the histidine residue in the middle contact position of the second ZF of Zif268 donates an H-bond to the N7 or O6 atom of guanine. The role of arginine and histidine residues in the interaction with guanine in ZF polypeptide-DNA complexes is confirmed by experiments of directed mutagenesis [5,6,9,14]. The crystallographic investigations of DNA-binding domains of lambda and phage 434 repressors, complexed with corresponding operator sites, revealed that guanine can also be H-bonded by lysine, asparagine, glutamine and serine residues [17,18]. No doubt, the remaining polar AA's -- threonine and tyrosine -- are able to form analogous bonds with guanine.

In fingers 1 and 3 of the Zif268-DNA complex, the second (middle) critical position is occupied by a glutamic acid that does not contact the cytosine at the corresponding region in the DNA [7]. However, ZF protein-DNA binding assays have shown that in natural binding sites this interaction does occur with both glutamic acid and aspartic acid [5,6,9,14,19]. Desjarlais and Berg [14] proposed an H-bonding formula for the interaction between cytosine and aspartic acid. The authors emphasized that the preference for aspartic or glutamic acid in the interaction with cytosine depends on the presence of glutamine or arginine in the third contact position ( $Z_3$ ), and serine or aspartic acid in the second position ( $Z_2$ ). The mutagenesis experiments of Nardelli et al. [5] reveal that cytosine can interact with a glutamine residue.



This may also be true for asparagine, which has similar polar groups. Cytosine should also be capable of making an H-bond with the hydroxyl oxygen atom in serine and threonine residues.

Thymine in the Zif268-DNA complex does not seem to participate in the recognition process. However, the crystal structure investigations of the lambda repressor, DNA-binding-domain DNA and engrailed homeodomain-DNA complexes, as well as ZF protein-DNA binding assays, demonstrate that thymine can make both hydrophobic contacts with non-polar residues (alanine, leucine, isoleucine, valine) and H-bonds with polar AA's (lysine, arginine, glutamine) [8,11,14,17,20].

The X-ray crystallographic studies of lambda and phage 434 repressor, DNA-binding domain complexes with corresponding operator sites revealed that an adenine base forms two H-bonds to glutamine: 1) the amide NH<sub>2</sub>-group of the glutamine side chain donates an H-bond to the N7-atom of adenine and 2) the amide O-atom accepts an H-bond from the N6 atom [17,18]. Similar H-bonds have been found between adenine and asparagine residues in the two homeodomain complexes [20,21]. ZF protein-DNA binding assays also indicate, that in ZF contact positions, adenine makes strong interactions with both glutamine and asparagine [8,11,12,14]. Considering that glutamic and aspartic acid carboxylic groups have O-atoms capable of accepting H-bonds as do glutamine and asparagine amide O-atoms, one may suppose that adenine can form a single H-bond with both glutamic and aspartic acid. Indeed, Letovsky and Dynan [19] have shown in a directed mutagenesis investigation that transcription factor Sp1, containing a glutamic acid residue in the central contact position of the ZF, binds only 3-fold more weakly to the adenine-substituted variant (-GAG-) than to the wild consensus recognition site (-GCG-). In addition, Desjarlais and Berg [14] and Berg [8] think it probable that adenine can (like guanine in the Zif268 -DNA complex) make one H-bond to a histidine residue. It is likely that not only histidine but also other polar amino acids (arginine, lysine, tyrosine, serine and threonine) are capable of forming an H-bond to atom N7 of adenine.

A database of potential ZF protein domains, containing 1,851 entries, has been assembled. This database was used computationally to observe the homologies between the ZF domains.

Several years ago Seeman et al. [22] concluded that a single H-bond is inadequate for uniquely identifying any particular base pair, as this leads to numerous degeneracies. They proposed that fidelity of recognition may be achieved using two H-bonds, as occurs in the major groove when asparagine or glutamine binds to adenine, and arginine binds to guanine.

On the basis of the above-given results, it was reasonable to test, using the models described herein, base recognition at the ZF contact positions of the following AA's:

- 1) guanine - R, H, K, Y, Q, N, S, T;
- 2) cytosine - E, D, Q, N, S, T;
- 3) thymine - I, L, V, A, R, H, K, Y, Q, N, S, T;
- 4) adenine - Q, N, E, D, H, R, K, Y, S, T

Plastic space-filling atomic-molecular and ionic models [23,24] have been used to build ZF-DNA complex imitations. These molecular models were chosen due to the extraordinary firmness of their connectors, their convenient scale ( $1\text{cm} = 1\text{Å} = 0.1\text{nm}$ ) and their improved theoretical parameters which were very suitable for the modeling of macromolecules. New modules of tetrahedral carbon atoms, with bond angles  $100^\circ$  and  $105^\circ$ , dihedral oxygen atoms ( $120^\circ$ ) and tetrahedral phosphorus atoms ( $102^\circ$  and  $118^\circ$ ), maintained the exact modeling of deoxyribose puckering and sugar-phosphate chain conformation in the B-DNA model. Peptide bonds in the DBP models were imitated by the fixing, to each other, of special modules of carbon atoms (bond angles  $116^\circ$ ,  $120.5^\circ$  and  $123.5^\circ$ ) and nitrogen atoms ( $122^\circ$  and  $119^\circ$ ). The zinc ion was represented in the model by a sphere ( $R = 0.85\text{ cm}$ ) fixed tetrahedrally to N and S atom modules of ZF histidine and cysteine residues. A long horizontal 34-base B-form DNA model with laterally-fixed DBP models was used for docking experiments.

In the first stage of the subject investigation, the models of Zif268 fingers 1, 2 and 3 were assembled, and the general spatial orientation of the ZF-B-DNA complex was observed. In the second stage, the steric fitness of all 64 nucleotide triplets to the different combinations of the above-mentioned AA's in the critical positions of the ZF-DNA complex was modeled.

A plastic molecular model of the Zif268 peptide-DNA complex was assembled on the basis of crystallographic data [7]. After the imitating of ZF-DNA backbone contacts and H-bonds between AA and bases in the major groove, it was confirmed that the overall arrangement of Zif268 is antiparallel to the DNA strand. The most steady ZF-DNA, nonspecific interaction seems to be the H-bond between a phosphodiester oxygen atom and the first invariant histidine residue fixed to the  $\text{Zn}^{2+}$  ion. A conserved arginine on the second b strand also contacts phosphodiester oxygen atoms on the primary DNA strand. However, fingers 2 and 3 of Zif268 contact equivalent phosphates with respect to the 3-bp sub-sites,

whereas the finger-1 H-bond is shifted by one nucleotide. Another four ZF-DNA backbone contacts made by arginine and serine residues are even more irregular in relation to the ZF modular structure.

All 11 critical H-bonds found in the Zif268-DNA crystal complex have been observed in the plastic models. As expected, the threonine residue in the first contact position of the second finger was too far from thymine to make an H-bond. However, differing from the results of crystal structure analysis, the model investigation clearly indicated the possibility of hydrogen bonding between a glutamic acid residue and cytosine in the second contact position of fingers 1 and 3.

It is noteworthy that, of the six guanine-AA contacts in recognition positions observed in the Zif268-DNA crystal structure, five were made with arginine and only one with histidine. It is even more interesting that this histidine-guanine interaction was the only one in the central-specific position. Considering the smaller size of histidine in comparison with arginine, it may be supposed that the middle position has steric constraints prohibiting contact between guanine and the larger arginine residue, although, due to its capability of forming two H-bonds, the latter pairing should be energetically favored.

To investigate the spatial conditions in different recognition positions, a B-DNA model was built which contained, in the primary strand, 1) the triplet GGG, and 2) models of ZF  $\alpha$ -helical protein fragments (including the AA immediately preceding the  $\alpha$ -helix) with a) side groups of the first Zn-binding histidine and b) groups for critical AA triplets  $R_1R_2R_3$  and  $R_1H_2R_3$ . The models of  $\alpha$ -helical fragments were fixed to the B-DNA model by an imitation of an H-bond joining a phosphodiester oxygen atom with a histidine residue. Specific base-AA contacts were then tested in these complexes. It was elucidated that only the complex GGG- $R_1H_2R_3$  contains the contact groups in positions corresponding to the distances of critical H-bonds found in the Zif268-DNA crystal structure. The complex GGG- $R_1R_2R_3$  is sterically unfavorable; molecular modeling reveals that, although in the outer contact positions guanine and arginine can be joined by two H-bonds, in the middle position such a pair cannot be included due to the limited space.

Observations derived from the physical models confirmed the supposition of steric constraints for some AA-base contacts in the central contact position. In the case of the complex  $G_1G_2G_3$ - $R_1H_2R_3$ , the following approximate distances from guanine N7 and O6 atoms to the C<sub>1</sub> atoms of corresponding AA's have been determined:  $G_1N7-R_1=7\text{\AA}$ ,  $G_1O6-R_1=8\text{\AA}$ ,  $G_2N7-H_2=5.5\text{\AA}$ ,  $G_2O6-H_2=6.5\text{\AA}$ ,  $G_3N7-R_3=8\text{\AA}$  and  $G_3O6-R_3=7\text{\AA}$ .

Using the models, the investigation of B-DNA and  $\alpha$ -helix basic structure elucidated the molecular basis for steric constraints in the second ZF-DNA recognition position. Joining, by a straight line, the analogous atomic groups (for example, N7 atoms of guanine) of the first and third base in the DNA triplet in the major groove results in the corresponding group of the middle (second) base being distanced from this line by about 1.5Å. Similarly, joining the C $\alpha$  atoms of the AA's in the first and third contact positions of the ZF by such a line results in the C $\alpha$  atom in the middle position also being at a distance of about 1.5Å. Thus, the space allowed for a critical AA in the middle contact position is compressed from both sides approximately 1.5Å.

Analysis of the above-given data on the ZF-DNA backbone contacts, as well as observations derived from the models, led to the conclusion that there are considerable differences in spatial conditions between first and third ZF-DNA recognition positions. In the first position the C $\alpha$  atom of the AA is distanced about 6.5Å from the phosphodiester oxygen atom where the ZF protein is fixed to the DNA backbone by the invariant histidine residue. Due to the steady fixing of this ZF  $\alpha$ -helical part by histidine, the freedom of conformational rearrangements in the first contact position is limited: the C $\alpha$  atom, with corresponding side chain, can be moved 2-3Å "up and down" in the plane of the base where it is localized in the primary DNA strand or, alternatively, 1-2Å perpendicularly to this plane.

On the other hand, the fixing of the N-terminal end of the ZF  $\alpha$ -helical region to the DNA backbone seems to be rather loose and variable, therefore allowing relatively large rearrangements for the C $\alpha$  atom and the corresponding AA in the third contact position. The latter contact position is favored by the fact that the C $\alpha$  atom in this position is more distant from the main fixation place (about 10.5Å from the phosphodiester atom bound to the histidine residue), and the corresponding AA in this position is not a part of the  $\alpha$ -helix. The most important finding is that, due to the above-described circumstances, the critical AA in the third contact site can apparently occupy very different positions in the corresponding bp plane. This means this residue may, in certain complexes, be very close to the base of the complementary DNA strand. One of the reasons for the appearance of such a geometrical configuration is that the typical, right-handed helical twist of B-DNA makes the complementary base on the nucleic acid second strand in the third contact site even more accessible than the main base on the primary chain. Molecular modeling clearly shows that in the third, and also partially in the second contact position, this DNA strand is capable of participating in the ZF-nucleic acid recognition process. In the Zif268-DNA crystal complex,

the  $\alpha$ -helix of each ZF domain, which is bound only to the DNA primary strand, is tipped at about a  $45^\circ$  angle with respect to the plane of the base pairs [7]. In cases wherein the second DNA strand, via critical H-bonds involving the third and second contact positions, is involved in the reading process, the direction of the  $\alpha$ -helix axis should be even more perpendicular to the base pair plane.

Thus, this more detailed investigation of ZF-DNA-complex imitations, through use of physical molecular models, shows that steric conditions in each of the three contact regions are different. These steric conditions are reflected in the ZF-DNA recognition rules.

On the basis of information obtained above, which yielded a general observation of steric conditions in the ZF-DNA recognition process, an extensive model study of various AA-base combinations in the critical contact positions was undertaken. The results of this investigation are presented both as the ZF-DNA reading code and main rules for recognition (Tables 1, 2 and 3). The rules are in good accordance with crystallographic, directed mutagenesis, DNA-binding and sequence analysis data.

With reference to the sequence of Formula I and the 2-dimensional structure diagram in Figure 2 (which provides a schematic representation of a zinc-finger domain and its interaction with a DNA strand), the studies confirmed the identity of the three critical contact positions in a given zinc-finger domain as follows:

- 1) between the first nucleotide in the triplet and the first AA preceding the constant histidine at the COOH end of the  $\alpha$ -helix;
- 2) between the second nucleotide in the triplet and the fourth AA preceding the constant histidine at the COOH end of the  $\alpha$ -helix; and,
- 3) between the third nucleotide in the triplet and the seventh AA preceding the constant histidine at the COOH end of the  $\alpha$ -helix.

Steric conditions in the three contact sites of the ZF-DNA recognition complexes are different. The first contact position is relatively large and strictly fixed, which enables the binding of a longer AA to bases on the primary DNA strand with sufficient specificity and affinity. The second position is compressed and can accommodate smaller AA's with somewhat lower specificity and affinity. The third position allows considerable conformational rearrangements including the contacts with the complementary DNA strand.

In Table 1, for each nucleotide of a given DNA triplet on the primary strand, both main (Column A) and alternative (Column B) base-binding AA's are presented. Both specificity and affinity were considered in including a residue in Column A. As was proposed already by

Seeman et al. [22], the fidelity of recognition is better maintained, in the case of purine bases (guanine and adenine), because they occupy a greater portion of the major groove and offer more hydrogen bonding sites than the pyrimidines. Therefore, the strongest AA interactions appeared to be those of arginine, glutamine and asparagine, each binding by two H-bonds to either guanine or adenine. The affinities of aspartic acid, glutamic acid, asparagine and glutamine were frequently enhanced by the formation of water bridges between carboxylate or amide oxygen atoms and DNA backbone, phosphodiester oxygen atoms. Although van der Waals interactions are relatively weak, they can play a certain role in recognition of the thymine methyl group by hydrophobic AA's (alanine, valine, leucine and isoleucine).

As indicated in Table 1, in many ZF-DNA complexes the base recognition in the nucleotide triplet of the primary DNA strand occurs not entirely via the primary strand, but by binding simultaneously to both the primary and complementary strands, or even exclusively to the complementary strand. Without "help" from the complementary DNA strand, the binding of critical AA's to nucleotides of the primary DNA chain would be too weak, in the case of several triplets, to realize the recognition process. All possible AA replacements were tested for strength of interaction in the  $Z_1$ - $Z_3$  positions. Domains with fewer than 2 hydrogen bonds on the primary strand were considered to be unstable.

Table 2 presents the ZF AA triplets having the highest affinity for interaction with corresponding DNA triplets. These ZF triplets contain only the main residues presented in Column A of Table 1. Table 2 also presents the binding energy components (H-bonds, water bridges, van der Waals interactions) maintaining the ZF-DNA recognition process in specific contact regions.

As can be seen from Table 2, the participation of the complementary DNA strand in the process of ZF binding, combined with the number of interactions (H-bonds, water bridges and van der Waals interactions) possible in the three contact regions, when optimal combinations are used, makes it possible to show that a complex formation with all 64 DNA triplets can be achieved. Table 2 shows that the maximal number of H-bonds, the strongest of the three types of interactions, is obtained when the first nucleotide of the triplet is guanine or adenine.

In nucleotide triplets wherein the number of H-bonds possible is less than maximal, the deficiency is often partially compensated by a significant amount of water-bridging between critical AA's and the sugar-phosphate backbone.

Even in cases wherein the first nucleotide of the triplet is thymine, and the number of the H-bonds is lowest, 1) the formation of two H-bonds between the AA in the Z<sub>3</sub> position, and the adenine and complementary thymine in the third contact position, and 2) probably, a single H-bond between thymine and serine or threonine in the second contact position, means that even TTN triplets can bind a ZF protein with sufficient affinity.

In any event, to obtain DBP's of the greatest effectiveness, attention should be paid to having the strongest interactions in the flanking contact points (1 and 3). If weaker combinations must be used, they would have less effect if positioned in the center contact point (2). It is important to note, however, that even weak binding in the contact points is important for establishing specificity.

Table 3 presents the main ZF AA triplets of Table 2, as well as the alternative AA's (shown in Column B of Table 1) which would be also expected to provide effective binding to the respective bases of a given DNA triplet. Table 3 also presents the binding energy components (H-bonds, water bridges, van der Waals interactions) maintaining the ZF-DNA recognition process in specific contact regions.

Table 1 - Z1

Codon	Z1 Column A	Z1 Column B	Hydrogen Bonds	Water Contacts	Hydrophobic Contacts
AAC	Q=	E*/R <sub>1</sub> -/K <sub>1</sub> -/N <sub>1</sub> =/D <sub>1</sub> */(H-/Y-/S-/T-)	6	0	0
AAG	Q=	E*/R-/K-	6	0	0
AAT	Q=	R-/K-/E*	6	0	0
ACC	Q=	E*/K <sub>1</sub> -	6	0	0
ACT	Q=	E-/R <sub>1</sub> -/K <sub>1</sub> -	6	0	0
GAA	R=	K-/H <sub>1</sub> -/Y <sub>1</sub> -/Q <sub>1</sub> -	6	0	0
GAC	R=	K-/H <sub>1</sub> -/Y <sub>1</sub> -	6	0	0
GAG	R=	H-/K-/Y-/Q-	6	0	0
GAT	R=	H-/K-/Y-/Q*	6	0	0
GCC	R=	H-/K-/Q-/N-/ (Y-/S-/T-)	6	0	0
GCT	R=	H-/K-/Y-/Q-	6	0	0
ACA	Q=	R-/K-/N-/E-/D-	5	1	0
ACG	Q=	R-/K-/N-/E*/D*	5	1	0
AGA	Q=	E*/R <sub>1</sub> -/K <sub>1</sub> -	5	1	0
AGG	Q=	E*/R <sub>1</sub> -/K <sub>1</sub> -	5	1	0
CAA	E*	Q*/N <sub>1</sub> */D <sub>1</sub> */R <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/S <sub>2</sub> -/T <sub>2</sub> -	5	1	0
CAG	E*	Q*/N*/D*/R <sub>2</sub> -/K <sub>2</sub> -	5	1	0
CAT	E*	Q*/R <sub>2</sub> -/K <sub>2</sub> -/(D*/N*/S-/T-/Y <sub>2</sub> -)	5	1	0
CCC	E*	Q*/N <sub>1</sub> */D <sub>1</sub> */R <sub>2</sub> =/K <sub>2</sub> -/Q <sub>2</sub> */N <sub>2</sub> *	5	1	0
CCT	E*	Q*/R <sub>2</sub> -/K <sub>2</sub> -	5	1	0
GCA	R=	K-/Q-/ (H-/Y-/N-/S-/T-)	5	1	0
GCG	R=	H-/K-/Y-/Q-/N-/S-/T-	5	1	0
GGA	R=	H-/K-/Q*/N*/Y <sub>1</sub> -	5	1	0
AAA	R-/K-		5	0	0
AGC	Q=	R-/K-/E*	5	0	0
AGT	Q=	E*/R <sub>1</sub> -/K <sub>1</sub> -	5	0	0
GGC	R=	K-/Q*/ (H-/Y-/N-)	5	0	0
GCG	R=	H-/K-/Y-/Q*/N*	5	0	0
GGT	R=	H-/K-/Y-/Q-/N-	5	0	0
CAC	E*	Q*/R <sub>2</sub> -/K <sub>2</sub> -	4	2	0
CCA	E*	Q*/R <sub>2</sub> -/K <sub>2</sub> -	4	2	0
CCG	E*	Q*/N*/D*/R <sub>2</sub> =/K <sub>2</sub> -	4	2	0
CGA	E*	Q*/N*/D*/R <sub>2</sub> -/K <sub>2</sub> -	4	2	0
CGC	E*	Q*/N*/D*/R <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> -/(S-/T-)	4	1	0
CGG	E*	Q*/N*/D*/R <sub>2</sub> -/K <sub>2</sub> -/Q <sub>2</sub> =	4	1	0
CGT	E*	Q*/D <sub>1</sub> */R <sub>2</sub> =/K <sub>2</sub> -/Q <sub>2</sub> -	4	1	0
ATA	Q=	E*/R <sub>1</sub> -/K <sub>1</sub> -	4	0	1
ATC	Q=	N-/E*/D*/R <sub>1</sub> -/K <sub>1</sub> -	4	0	1
ATG	Q=	R-/K-/E-/ (H-)	4	0	1
ATT	Q=	E*/R <sub>1</sub> -/K <sub>1</sub> -	4	0	1
GTA	R=	H-/K-/Y-/Q-	4	0	1
GTC	R=	H-/K-/Y-/Q*	4	0	1
GTG	R=	K-/Q-/H <sub>1</sub> -	4	0	1
GTT	R=	H-/K-/Q-/N-/Y <sub>1</sub> -	4	0	1
TAA	1#/L#	R-/K-/Q-/V <sub>1</sub> #	4	0	1
TAG	1#/L#	R-/K-/Q*	4	0	1
TAT	1#/L#/V#	R-/K-/Y-/Q-/N-	4	0	1
TCC	1#/L#	R-/H-/K-/Q*/N*/V#A#	4	0	1
TCT	1#/L <sub>1</sub> #	R-/H-/K-/Q*	4	0	1



Table 1 - Z1

CTA	E*	Q*	3	1	1
CTC	E*	Q*/N <sub>1</sub> -/D <sub>1</sub> -/R <sub>2</sub> =/K <sub>2</sub> -/Q <sub>2</sub> -/N <sub>2</sub> -	3	1	1
CTG	E*	Q*/N*/D*/R <sub>2</sub> -/K <sub>2</sub> -	3	1	1
CTT	E*	Q*/N*/D*/R <sub>2</sub> -/K <sub>2</sub> -/Q <sub>2</sub> =	3	1	1
TCA	1#/L#	R-/K-/Q*	3	1	1
TCG	1#/L#	R-/K-/Q*	3	1	1
TGA	1#/L#	R-/K-/Q*	3	1	1
TAC	1#/L#/V#	R-/K-/Q*	3	0	1
TGC	1 <sub>1</sub> #/L <sub>1</sub> #/V <sub>1</sub> #	R-/K-/H <sub>1</sub> -/Q <sub>1</sub> */N <sub>1</sub> */S <sub>1</sub> -/T <sub>1</sub> -	3	0	1
TGG	1#/L#	R-/K-/Q*	3	0	1
TGT	1#	R-/H-/K-/Q*/N*/L#	3	0	1
TTA	1#/L#	R-/K-/Q*/N*	2	0	2
TTC	1#/L#/V#	R-/K-/Q*/N*	2	0	2
TTG	1#/L#	R-/K-/Y-/Q*	2	0	2
TTT	1#/L#	R-/K-/Q*/N*/V <sub>1</sub> #	2	0	2

## Legend

where / separates alternative amino acids

where X without subscript has all its interactions with the primary strand

where  $X_1$  has some interactions with the primary strand and some interactions with the complementary strand

where  $X_2$  has interaction with the complementary strand

where  $X_3$  has interactions with both the primary and complementary strands

where - is one hydrogen bond between the amino acid and the base

where = is two hydrogen bonds between the amino acid and the base

where \* is one hydrogen bond via a water bridge between the amino acid and the phosphodiester oxygen atom of the backbone

where # is one or more van der Waals contacts between the amino acid and the base

where amino acids in ( ) have interaction with the base of the primary strand where one of two other possible protein-DNA recognition interactions is absent

Table 1 - Z2

Codon	Z2 Column A	Z2 Column B	Hydrogen Bonds	Water Contacts	Hydrophobic Contacts
AAC	Q <sub>1</sub> =	N=/D*/S-/T-/R <sub>1</sub> -/K <sub>1</sub> -/E <sub>1</sub> */(H-/Y-)	6	0	0
AAG	Q=/N=	R-/H-/K-/E*/D*/K <sub>2</sub> -	6	0	0
AAT	N=	K-/Q=/E*/D*/R <sub>1</sub> -/H <sub>1</sub> -/K <sub>2</sub> -/Q <sub>2</sub> =/N <sub>2</sub> =	6	0	0
ACC	Q <sub>2</sub> =/E*	D*/S-/T-/N <sub>2</sub> =/(K <sub>2</sub> -)	6	0	0
ACT	N <sub>2</sub> =/D*	Q*/N*/E*/S-/T-/K <sub>2</sub> -/Q <sub>2</sub> =	6	0	0
GAA	N=	D*/R <sub>1</sub> -/H <sub>1</sub> -/K <sub>1</sub> -/Y <sub>1</sub> -/Q <sub>1</sub> =/E <sub>1</sub> */K <sub>2</sub> -	6	0	0
GAC	N=	D*/R <sub>1</sub> -/H <sub>1</sub> -/K <sub>1</sub> -/Y <sub>1</sub> -/Q <sub>1</sub> =/E <sub>1</sub> */K <sub>2</sub> -	6	0	0
GAG	N=	Q=/E*/D*/R <sub>1</sub> -/H <sub>1</sub> -/K <sub>1</sub> -/Y <sub>1</sub> -/K <sub>2</sub> -	6	0	0
GAT	Q=/N=	K-/E*/D*/K <sub>2</sub> -/(R-/H-/Y-)	6	0	0
GCC	Q <sub>2</sub> =/E*	Q*/N*/D*/S-/T-/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Q <sub>2</sub> =/N <sub>2</sub> =/S <sub>2</sub> -/T <sub>2</sub> -	6	0	0
GCT	N <sub>2</sub> =/D*	Q-/N-/E-/S-/T-/H <sub>2</sub> -/K <sub>2</sub> -/N <sub>2</sub> */Q <sub>2</sub> =/(R <sub>2</sub> =/Y <sub>2</sub> -)	6	0	0
ACA	D*	Q*/N*/E*/S-/T-/K <sub>2</sub> -	5	1	0
ACG	E*/D*	Q*/N*/S-/T-/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> =/N <sub>2</sub> =	5	1	0
AGA	N*	R <sub>1</sub> -/H <sub>1</sub> -/K <sub>1</sub> -/Y <sub>1</sub> -/Q <sub>1</sub> *	5	1	0
AGG	Q*/N*	R <sub>1</sub> */H <sub>1</sub> */K <sub>1</sub> -	5	1	0
CAA	N=	D*/S-/T-/R <sub>1</sub> -/H <sub>1</sub> -/K <sub>1</sub> -/Y <sub>1</sub> -/Q <sub>1</sub> =/E <sub>1</sub> -	5	1	0
CAG	Q=	N=/E*/D*/R <sub>1</sub> -/K <sub>1</sub> -/K <sub>2</sub> -/Q <sub>2</sub> =	5	1	0
CAT	Q <sub>1</sub> =/N=	D-/S-/T-/R <sub>1</sub> -/H <sub>1</sub> -/K <sub>1</sub> -/Y <sub>1</sub> -/E <sub>1</sub> =/K <sub>2</sub> -/Q <sub>2</sub> =/N <sub>2</sub> =	5	1	0
CCC	Q <sub>2</sub> =/E <sub>1</sub> *	N*/D*/S-/T-/H <sub>2</sub> -/K <sub>2</sub> -/N <sub>2</sub> =/(Y <sub>2</sub> -)	5	1	0
CCT	N <sub>2</sub> =/D*	Q*/N*/E*/S-/T-/H <sub>2</sub> -/K <sub>2</sub> -/Q <sub>2</sub> =	5	1	0
GCA	D*	Q*/N*/E*/S-/T-/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> =/N <sub>2</sub> =	5	1	0
GCG	E*/D*	Q*/N*/S-/T-/K <sub>2</sub> -/Q <sub>2</sub> =/N <sub>2</sub> =/(H <sub>2</sub> -)	5	1	0
GGA	N*	Q*/S-/T-/K <sub>1</sub> -/(R-/Y-/H-)	5	1	0
AAA	N=	D*/R <sub>1</sub> -/H <sub>1</sub> -/K <sub>1</sub> -/Y <sub>1</sub> -/Q <sub>1</sub> =/E <sub>1</sub> *	5	0	0
AGC	H <sub>1</sub> -	Q*/N*/S-/T-/R <sub>1</sub> -/K <sub>1</sub> -/(R=/Y-)	5	0	0
AGT	H <sub>1</sub> -	Q*/N*/S-/T-/R <sub>1</sub> =/K <sub>1</sub> -	5	0	0
GGC	H <sub>1</sub> -	N*/S-/T-/K <sub>1</sub> -/(R-/K-/Y-/Q*)	5	0	0
GGG	H-	K-/Q*/N*/S-/T-/Y <sub>1</sub> -/(R-)	5	0	0
GGT	H-	Q*/N*/S-/T-/R <sub>1</sub> -/K <sub>1</sub> -	5	0	0
CAC	N=	D*/R <sub>1</sub> -/K <sub>1</sub> -/Q <sub>1</sub> =/E <sub>1</sub> *	4	2	0
CCA	D*	N*/S-/T-/Q <sub>1</sub> */E <sub>1</sub> */K <sub>2</sub> -/N <sub>2</sub> =	4	2	0
CCG	D*	Q*/N*/E*/S-/T-/K <sub>2</sub> -/Q <sub>2</sub> =/N <sub>2</sub> =/(R-/H-/Y-)	4	2	0
CGA	N*	Q*/S-/T-/R <sub>1</sub> =/H <sub>1</sub> -/K <sub>1</sub> -/(Y-)	4	2	0
CGC	H <sub>1</sub> -	Q*/N*/S-/T-/K <sub>1</sub> -/(R <sub>1</sub> =/Y-)	4	1	0
CGG	H-	K-/Q*/N*/R <sub>1</sub> -/(Y-)	4	1	0
CGT	H-	Q*/N*/R <sub>1</sub> -/K <sub>1</sub> -/(Y-)	4	1	0
ATA	1#/L#/V#/A#	S-/T-/K <sub>1</sub> -/Q <sub>1</sub> */N <sub>1</sub> */(R-/H-/Y-)	4	0	1
ATC	1#/L#	N*/S-/T-/V-/R <sub>1</sub> -/H <sub>1</sub> -/K <sub>1</sub> -/Y <sub>1</sub> -/Q <sub>1</sub> */R <sub>2</sub> -/K <sub>2</sub> -/Q <sub>2</sub> =/N <sub>2</sub> =/E <sub>2</sub> */D <sub>2</sub> *	4	0	1
ATG	1#/L#/V#	S-/T-/E <sub>2</sub> -/D <sub>2</sub> -/Q <sub>2</sub> =/N <sub>2</sub> =/(R-/H-/K-/Y-)	4	0	1
ATT	1#/L#/V#	R-/H-/K-/Y-/Q*/N*/S-/T-	4	0	1
GTA	1#/L#/V#/A#	Q*/N*/S-/T-/R <sub>1</sub> -/H <sub>1</sub> -/K <sub>1</sub> -/E <sub>2</sub> */D <sub>2</sub> */Q <sub>2</sub> =/N <sub>2</sub> =/(Y-)	4	0	1
GTC	1#/L#	N-/S-/T-/V-/R <sub>1</sub> -/H <sub>1</sub> -/K <sub>1</sub> -/Q <sub>1</sub> *	4	0	1
GTG	1#/L#/V#	Q*/N*/S-/T-/H <sub>1</sub> -/K <sub>1</sub> -	4	0	1
GTT	1#/L#/V#	Q*/N*/S-/T-/K <sub>1</sub> -/(R-/H-/Y-/A-)	4	0	1
TAA	N=	D*/R <sub>1</sub> -/H <sub>1</sub> -/K <sub>1</sub> -/Y <sub>1</sub> -/Q <sub>1</sub> =/E <sub>1</sub> -	4	0	1

Table 1 - Z2

TAG	N=	Q=/E*/D*/R <sub>1</sub> -/K <sub>1</sub> -/K <sub>2</sub> -	4	0	1
TAT	N=	K-/Q=/N=/E*/D*/S-/T-/H <sub>1</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Q <sub>2</sub> =/N <sub>2</sub> =/(R-/Y-)	4	0	1
TCC	Q <sub>2</sub> =/E*	Q*/N*/D*/S-/T-/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> -/N <sub>2</sub> -	4	0	1
TCT	N <sub>2</sub> =/D*	Q*/N*/E*/S-/T-/R <sub>2</sub> =/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> *	4	0	1
CTA	1#/L#/V#/A#	S-/T-/Q <sub>1</sub> =/N <sub>1</sub> =/(H-/K-)	3	1	1
CTC	1#/L#	S-/T-/V-/R <sub>1</sub> -/H <sub>1</sub> -/K <sub>1</sub> -/Y <sub>1</sub> -/Q <sub>1</sub> */N <sub>1</sub> *	3	1	1
CTG	1#/L#/V#	N*/S-/T-/K <sub>1</sub> -/Q <sub>1</sub> *	3	1	1
CTT	1#/L#	Q*/N*/S-/T-/V#/R <sub>1</sub> -/H <sub>1</sub> -/K <sub>1</sub> -/E <sub>2</sub> -/D <sub>2</sub> -/Q <sub>2</sub> =/N <sub>2</sub> =/(Y-)	3	1	1
TCA	D*	Q*/N*/E*/S-/T-/R <sub>2</sub> =/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> =/N <sub>2</sub> =	3	1	1
TGG	D*	Q*/N*/E*/S-/T-/R <sub>2</sub> =/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/S <sub>2</sub> -/T <sub>2</sub> -/Q <sub>2</sub> =/N <sub>2</sub> =	3	1	1
TGA	N*	Q*/S-/T-/H <sub>1</sub> -/K <sub>1</sub> -/(R-/Y-)	3	1	1
TAC	N=	D-/H <sub>1</sub> -/K <sub>1</sub> -/Q <sub>1</sub> =/E <sub>1</sub> */K <sub>2</sub> -/(R-/Y-)	3	0	1
TGC	H <sub>1</sub> -	Q*/N*/S-/T-/R <sub>1</sub> =/K <sub>1</sub> -/Y <sub>1</sub> -	3	0	1
TGG	H-	R-/K-/Q*/N*/Y <sub>1</sub> -	3	0	1
TGT	H <sub>1</sub> -	N*/S-/T-/K <sub>1</sub> -/Y <sub>1</sub> -/Q <sub>1</sub> */(R=)	3	0	1
TTA	1#/L#/V#/A#	S-/T-/R <sub>1</sub> -/H <sub>1</sub> -/K <sub>1</sub> -/Y <sub>1</sub> -/E <sub>2</sub> -/D <sub>2</sub> -/Q <sub>2</sub> =/N <sub>2</sub> =	2	0	2
TTC	1#/L#	N*/S-/T-/V-/A-/R <sub>1</sub> -/H <sub>1</sub> -/K <sub>1</sub> -/Y <sub>1</sub> -/K <sub>2</sub> -/Q <sub>2</sub> =/N <sub>2</sub> =	2	0	2
TTG	1#/L#/V#	N*/S-/T-/H <sub>1</sub> -/K <sub>1</sub> -/Q <sub>1</sub> *	2	0	2
TTT	1#/L#/V#	Q*/N*/S-/T-/A-/H <sub>1</sub> -/K <sub>1</sub> -/(R-/Y-)	2	0	2

## Legend

where / separates alternative amino acids

where X without subscript has all its interactions with the primary strand

where  $X_1$  has some interactions with the primary strand and some interactions with the complementary strand

where  $X_2$  has interaction with the complementary strand

where  $X_3$  has interactions with both the primary and complementary strands

where - is one hydrogen bond between the amino acid and the base

where = is two hydrogen bonds between the amino acid and the base

where \* is one hydrogen bond via a water bridge between the amino acid and the phosphodiester oxygen atom of the backbone

where # is one or more van der Waals contacts between the amino acid and the base

where amino acids in ( ) have interaction with the base of the primary strand where one of two other possible protein-DNA recognition interactions is absent

Table 1 - Z3

Codon	Z3 Column A	Z3 Column B	Hydrogen Bonds	Water Contacts	Hydrophobic Contacts
AAC	R <sub>2</sub> =	Y-/Q*/N*/E*/D*/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> */N <sub>2</sub> *	6	0	0
AAG	R=	H-/K-/Y-/Q-/Q <sub>2</sub> */E <sub>2</sub> */N <sub>2</sub> */D <sub>2</sub> */S <sub>2</sub> -/T <sub>2</sub> -	6	0	0
AAT	Q <sub>2</sub> =	R-/H-/K-/Y-/Q*/R <sub>2</sub> /K <sub>2</sub> /Q <sub>2</sub> /N <sub>2</sub> /E <sub>2</sub> */D <sub>2</sub> */Q <sub>2</sub> =	6	0	0
ACC	R <sub>2</sub> =	Q*/E*/K <sub>2</sub> */N <sub>2</sub> */S <sub>2</sub> -/T <sub>2</sub> -	6	0	0
ACT	Q <sub>2</sub> =	R-/H-/K-/Y-/Q*/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/N <sub>2</sub> =/E <sub>2</sub> */D <sub>2</sub> *	6	0	0
GAA	Q=	R-/H-/K-/E*/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> *	6	0	0
GAC	Q <sub>2</sub> =/E*	Q*/R <sub>2</sub> =/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> */N <sub>2</sub> */S <sub>2</sub> -/T <sub>2</sub> -	6	0	0
GAG	R=	K-/Q*/Q <sub>2</sub> */E <sub>2</sub> */N <sub>2</sub> -/D <sub>2</sub> -/S <sub>2</sub> -/T <sub>2</sub> -	6	0	0
QAT	Q <sub>2</sub> =/Q <sub>2</sub> =	R-/H-/K-/Y-/Q*/N*/I-/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/N <sub>2</sub> =/E <sub>2</sub> */D <sub>2</sub> *	6	0	0
CCC	R <sub>2</sub> =	Q*/N*/E*/D*/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> */N <sub>2</sub> */S <sub>2</sub> -/T <sub>2</sub> -/Q <sub>2</sub> =/N <sub>2</sub> =	6	0	0
GCT	Q <sub>2</sub> =	R-/K-/Y-/Q/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/N <sub>2</sub> */E <sub>2</sub> */D <sub>2</sub> */S <sub>2</sub> -/T <sub>2</sub> -	6	0	0
ACA	Q=	R-/H-/K-/Y-/N=I/E*/D*/I#L#V#R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> */N <sub>2</sub> */(S <sub>2</sub> -/T <sub>2</sub> -)	5	1	0
ACG	R=	H-/K-/Y-/Q*/N*/S-/T-/E <sub>2</sub> */D <sub>2</sub> */Q <sub>2</sub> =/N <sub>2</sub> =	5	1	0
AGA	Q=	E-/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> -/N <sub>2</sub> */I <sub>2</sub> #L <sub>2</sub> #V <sub>2</sub> #A <sub>2</sub> #	5	1	0
AGG	R=	K-/Y-/Q <sub>2</sub> */N <sub>2</sub> */E <sub>2</sub> */D <sub>2</sub> *	5	1	0
CAA	Q <sub>2</sub> =	R-/H-/K-/Y-/Q=I/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> *	5	1	0
CAG	R=	H-/K-/Y-/Q-/N-/Q <sub>2</sub> */N <sub>2</sub> */E <sub>2</sub> */D <sub>2</sub> */S <sub>2</sub> -/T <sub>2</sub> -	5	1	0
CAT	Q <sub>2</sub> =/N <sub>2</sub> =	R-/H-/K-/Y-/Q*/D <sub>2</sub> -/E <sub>2</sub> -	5	1	0
CCC	R <sub>2</sub> =	Q*/E*/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> */N <sub>2</sub> */S <sub>2</sub> -/T <sub>2</sub> -	5	1	0
CCT	Q <sub>2</sub> =	R-/H-/K-/Y-/Q*/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/N <sub>2</sub> */E <sub>2</sub> */S <sub>2</sub> -/T <sub>2</sub> -	5	1	0
GCA	Q=Q <sub>2</sub> =	R-/H-/K-/Y-/E*/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> */N <sub>2</sub> */S <sub>2</sub> -/T <sub>2</sub> -/I <sub>2</sub> #L <sub>2</sub> #V <sub>2</sub> #A <sub>2</sub> #	5	1	0
GCG	R=	H-/K-/Y-/Q-/E <sub>2</sub> -/D <sub>2</sub> -/S <sub>2</sub> -/T <sub>2</sub> -/Q <sub>2</sub> =/N <sub>2</sub> =	5	1	0
GGA	Q=	N=I/E*/D*/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> */N <sub>2</sub> */I <sub>2</sub> #L <sub>2</sub> #V <sub>2</sub> #A <sub>2</sub> #	5	1	0
AAA	Q=	R-/K-/N=I/E-/D-/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> -/N <sub>2</sub> -/I <sub>2</sub> #L <sub>2</sub> #	5	0	0
AGC	R <sub>2</sub> =	Q*/N*/E*/D*/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> */N <sub>2</sub> *	5	0	0
AGT	Q <sub>2</sub> =	R-/K-/Y-/Q*/N*/S-/T-/I#L#H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/N <sub>2</sub> =I/E <sub>2</sub> */D <sub>2</sub> *	5	0	0
GGC	R <sub>2</sub> =	Q*/N*/E*/D*/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> */N <sub>2</sub> */S <sub>2</sub> -/T <sub>2</sub> -	5	0	0
GCG	R=	H-/K-/Y-/Q*/N*/E <sub>2</sub> */D <sub>2</sub> */S <sub>2</sub> -/T <sub>2</sub> -/Q <sub>2</sub> =/N <sub>2</sub> =	5	0	0
GGT	Q <sub>2</sub> =/N <sub>2</sub> =	R-/K-/Q*/N*/K <sub>2</sub> -/E <sub>2</sub> */D <sub>2</sub> */S <sub>2</sub> -/T <sub>2</sub> -	5	0	0
CAC	E*	Q*/R <sub>2</sub> =/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> *	4	2	0
CCA	Q=	E*/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> */N <sub>2</sub> */S <sub>2</sub> -/T <sub>2</sub> -/I <sub>2</sub> #L <sub>2</sub> #V <sub>2</sub> #A <sub>2</sub> #	4	2	0
CCG	R=	H-/K-/Q-/N-/E <sub>2</sub> */D <sub>2</sub> */Q <sub>2</sub> =/N <sub>2</sub> =	4	2	0
GGA	Q=	R-/H-/K-/Y-/N=I/E*/D*/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> */(N <sub>2</sub> */S <sub>2</sub> -/T <sub>2</sub> -/I <sub>2</sub> #L <sub>2</sub> #V <sub>2</sub> #A <sub>2</sub> #)	4	2	0
GGC	R <sub>2</sub> =	Q*/N*/E*/D*/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> -/N <sub>2</sub> -	4	1	0
GGG	R=	H-/K-/Y-/Q*/E <sub>2</sub> -/D <sub>2</sub> -/S <sub>2</sub> -/T <sub>2</sub> -/Q <sub>2</sub> =/N <sub>2</sub> =	4	1	0
GGT	Q <sub>2</sub> =/N <sub>2</sub> =	R-/K-/Y-/Q*/N*/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/E <sub>2</sub> */D <sub>2</sub> */Q <sub>2</sub> =	4	1	0
ATA	Q=	R-/K-/E*/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/I <sub>2</sub> #L <sub>2</sub> #V <sub>2</sub> #Q <sub>2</sub> =	4	0	1
ATC	R <sub>2</sub> =/E*	Q*/R <sub>2</sub> =/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> */N <sub>2</sub> *	4	0	1
ATG	R=	H-/K-/Y-/Q-/Q <sub>2</sub> */N <sub>2</sub> */E <sub>2</sub> */D <sub>2</sub> *	4	0	1
ATT	Q <sub>2</sub> =	R-/H-/K-/Y-/Q*/N*/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> =/N <sub>2</sub> =/E <sub>2</sub> */D <sub>2</sub> */S <sub>2</sub> -/T <sub>2</sub> -	4	0	1
GTA	Q=	R-/K-/Y-/N=I/E*/D-/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> -/N <sub>2</sub> -/S <sub>2</sub> -/T <sub>2</sub> -/I <sub>2</sub> #L <sub>2</sub> #V <sub>2</sub> #A <sub>2</sub> #	4	0	1
GTC	R <sub>2</sub> =/E*	Q*/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> */N <sub>2</sub> */S <sub>2</sub> -/T <sub>2</sub> -	4	0	1
GTG	R=	H-/K-/Y-/Q*/Q <sub>2</sub> */N <sub>2</sub> */E <sub>2</sub> */D <sub>2</sub> */S <sub>2</sub> -/T <sub>2</sub> -	4	0	1
GTT	Q <sub>2</sub> =	R-/H-/K-/Y-/Q*/N*/I#L#V#R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> =/N <sub>2</sub> =/E <sub>2</sub> */D <sub>2</sub> -	4	0	1
TAA	Q=	R-/H-/K-/Y-/E*/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> */I <sub>2</sub> #L <sub>2</sub> #V <sub>2</sub> #A <sub>2</sub> #/Q <sub>2</sub> =	4	0	1
TAG	R=	H-/K-/Y-/Q*/N*/Q <sub>2</sub> */N <sub>2</sub> -/E <sub>2</sub> */D <sub>2</sub> -	4	0	1
TAT	Q <sub>2</sub> =	R-/H-/K-/Y-/Q*/N*/S-/T-/I#L#V#R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> =/N <sub>2</sub> =/E <sub>2</sub> */D <sub>2</sub> */S <sub>2</sub> -/T <sub>2</sub> -	4	0	1
TCC	R <sub>2</sub> =	Q*/E*/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> -/N <sub>2</sub> -/S <sub>2</sub> -/T <sub>2</sub> -	4	0	1
TCT	Q <sub>2</sub> =/N <sub>2</sub> =	R-/H-/K-/Y-/Q*/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/E <sub>2</sub> -/D <sub>2</sub> -	4	0	1
CTA	Q=	R-/K-/E-/R <sub>2</sub> -/K <sub>2</sub> -/Q <sub>2</sub> */N <sub>2</sub> */S <sub>2</sub> -/T <sub>2</sub> -/I <sub>2</sub> #L <sub>2</sub> #V <sub>2</sub> #A <sub>2</sub> #	3	1	1
CTC	R <sub>2</sub> =/E*	Q*/N*/D*/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> -	3	1	1

Table 1 - Z3

CTG	R=	K-/Y-/Q <sub>2</sub> -/N <sub>2</sub> */E <sub>1</sub> -/D <sub>2</sub> *	3	1	1
CTT	Q <sub>2</sub> =	R-/H-/K-/Q*/N*/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/N <sub>2</sub> =/E <sub>2</sub> */D <sub>2</sub> *	3	1	1
TCA	Q=	Y-/E*/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/N <sub>2</sub> */S <sub>2</sub> -/T <sub>2</sub> -/L <sub>2</sub> */V <sub>2</sub> */A <sub>2</sub> */Q <sub>3</sub> =/N <sub>3</sub> =	3	1	1
TCG	R=	H-/K-/Y-/Q*/Q <sub>2</sub> */N <sub>2</sub> -/E <sub>2</sub> */D <sub>2</sub> -/S <sub>2</sub> -/T <sub>2</sub> -	3	1	1
TGA	Q=	R-/H-/K-/Y-/N= /E*/D*/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> -/Q <sub>3</sub> =	3	1	1
TAC	E-	Q-/R <sub>2</sub> =/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> */N <sub>2</sub> *	3	0	1
TGC	R <sub>2</sub> =	Q*/N*/E*/D*/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> */N <sub>2</sub> -	3	0	1
TGG	R=	H-/K-/Y-/Q-/N <sub>2</sub> -/E <sub>2</sub> -/D <sub>2</sub> -/S <sub>2</sub> -/T <sub>2</sub> -	3	0	1
TGT	Q <sub>2</sub> =	R-/H-/K-/Y-/Q*/L#/L#/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/N <sub>2</sub> =/E <sub>2</sub> *	3	0	1
TTA	Q=	R-/H-/K-/Y-/E*/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/N <sub>2</sub> -/S <sub>2</sub> -/T <sub>2</sub> -/L <sub>2</sub> */V <sub>2</sub> */A <sub>2</sub> */	2	0	2
TTC	R <sub>2</sub> =/E*	Q*/N*/D*/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> -/N <sub>2</sub> -/S <sub>2</sub> -/T <sub>2</sub> -	2	0	2
TTG	R=	H-/K-/Y-/Q-/Q <sub>2</sub> */N <sub>2</sub> */E <sub>2</sub> */D <sub>2</sub> *	2	0	2
TTT	Q <sub>3</sub> =	R-/H-/K-/Y-/R <sub>2</sub> -/H <sub>2</sub> -/K <sub>2</sub> -/Y <sub>2</sub> -/Q <sub>2</sub> =/N <sub>2</sub> =/E <sub>2</sub> -/D <sub>2</sub> -/N <sub>3</sub> =	2	0	2

## Legend

where / separates alternative amino acids

where X without subscript has all its interactions with the primary strand

where X<sub>1</sub> has some interactions with the primary strand and some interactions with the complementary strand

where X<sub>2</sub> has interaction with the complementary strand

where X<sub>3</sub> has interactions with both the primary and complementary strands

where - is one hydrogen bond between the amino acid and the base

where = is two hydrogen bonds between the amino acid and the base

where \* is one hydrogen bond via a water bridge between the amino acid and the phosphodiester oxygen atom of the backbone

where # is one or more van der Waals contacts between the amino acid and the base

where amino acids in ( ) have interaction with the base of the primary strand where one of two other possible protein-DNA recognition interactions is absent



Table 2

Codon	Z1 Column A	Z2 Column A	Z3 Column A	Hydrogen Bonds	Water Contacts	Hydrophobic Contacts
AAC	Q	Q	R	6	0	0
AAG	Q	N/Q	R	6	0	0
AAT	Q	N	Q	6	0	0
ACC	Q	E/Q	R	6	0	0
ACT	Q	D/N	Q	6	0	0
GAA	R	N	Q	6	0	0
GAC	R	N	E/Q	6	0	0
GAG	R	N	R	6	0	0
GAT	R	N/Q	Q	6	0	0
GCC	R	E/Q	R	6	0	0
GCT	R	D/N	Q	6	0	0
ACA	Q	D	Q	5	1	0
ACG	Q	D/E	R	5	1	0
AGA	Q	N	Q	5	1	0
AGG	Q	N/Q	R	5	1	0
CAA	E	N	Q	5	1	0
CAG	E	Q	R	5	1	0
CAT	E	N/Q	N/Q	5	1	0
CCC	E	E/Q	R	5	1	0
CCT	E	D/N	Q	5	1	0
GCA	R	D	Q	5	1	0
GCG	R	D/E	R	5	1	0
GGA	R	N	Q	5	1	0
AAA	K/R	N	Q	5	0	0
AGC	Q	H	R	5	0	0
AGT	Q	H	Q	5	0	0
GGC	R	H	R	5	0	0
GGG	R	H	R	5	0	0
GGT	R	H	N/Q	5	0	0
CAC	E	N	E	4	2	0
CCA	E	D	Q	4	2	0
CCG	E	D	R	4	2	0
CGA	E	N	Q	4	2	0
CGC	E	H	R	4	1	0
CGG	E	H	R	4	1	0
CGT	E	H	N/Q	4	1	0
ATA	Q	A/I/L/V	Q	4	0	1
ATC	Q	I/L	E/R	4	0	1

Table 2

ATG	Q	I/L/V	R	4	0	1
ATT	Q	I/L/V	Q	4	0	1
GTA	R	A/I/L/V	Q	4	0	1
GTC	R	I/L	E/R	4	0	1
GTG	R	I/L/V	R	4	0	1
GTT	R	I/L/V	Q	4	0	1
TAA	I/L	N	Q	4	0	1
TAG	I/L	N	R	4	0	1
TAT	I/L/V	N	Q	4	0	1
TCC	I/L	E/Q	R	4	0	1
TCT	I/L	D/N	N/Q	4	0	1
CTA	E	A/I/L/V	Q	3	1	1
CTC	E	I/L	E/R	3	1	1
CTG	E	I/L/V	R	3	1	1
CTT	E	I/L	Q	3	1	1
TCA	I/L	D	Q	3	1	1
TCG	I/L	D	R	3	1	1
TGA	I/L	N	Q	3	1	1
TAC	I/L/V	N	E	3	0	1
TGC	I/L/V	H	R	3	0	1
TGG	I/L	H	R	3	0	1
TGT	I	H	Q	3	0	1
TTA	I/L	A/I/L/V	Q	2	0	2
TTC	I/L/V	I/L	E/R	2	0	2
TTG	I/L	I/L/V	R	2	0	2
TTT	I/L	I/L/V	Q	2	0	2

where / separates alternative amino acids

Table

Codon	Z1 Column A	Z1 Column B	Z2 Column A	Z2 Column B	Z3 Column A	Z3 Column B	Hydrogen Bonds	Water Contacts	Hydrophobic Contacts
AAC	Q	D/E/H/K/N/R/S/T/Y	Q	D/E/H/K/N/R/S/T/Y	R	D/E/H/K/N/Q/Y	6	0	0
AAG	Q	E/K/R	N/Q	D/E/H/K/R	R	D/E/H/K/N/Q/S/T/Y	6	0	0
AAT	Q	E/K/R	N	D/E/H/K/N/Q/R	Q	D/E/H/K/N/Q/R/Y	6	0	0
ACC	Q	E/K	E/Q	D/K/N/S/T	R	E/K/N/Q/S/T	6	0	0
ACT	Q	E/K/R	D/N	E/K/N/Q/S/T	Q	D/E/H/K/N/Q/R/Y	6	0	0
GAA	R	H/K/Q/Y	N	D/E/H/K/Q/R/Y	Q	E/H/K/Q/R/Y	6	0	0
GAG	R	H/K/Y	N	D/E/H/K/Q/R/Y	E/Q	H/K/N/Q/R/S/T/Y	6	0	0
GAC	R	H/K/Q/Y	N	D/E/H/K/Q/R/Y	R	D/E/K/N/Q/S/T	6	0	0
GAT	R	H/K/Q/Y	N/Q	D/E/H/K/R/Y	Q	D/E/H/K/N/Q/R/Y	6	0	0
GCC	R	H/K/N/Q/S/T/Y	E/Q	D/H/K/N/Q/R/S/T	R	D/E/H/K/N/Q/S/T/Y	6	0	0
GCT	R	H/K/Q/Y	D/N	E/H/K/N/Q/R/S/T/Y	Q	D/E/H/K/N/Q/R/S/T/Y	6	0	0
ACA	Q	D/E/K/N/R	D	E/K/N/Q/S/T	Q	D/E/H/K/L/N/Q/R/S/T/V/Y	5	1	0
ACG	Q	D/E/K/N/R	D/E	H/K/N/Q/R/S/T/Y	R	D/E/H/K/N/Q/S/T/Y	5	1	0
AGA	Q	E/K/R	N	H/K/Q/R/Y	Q	A/E/H/K/L/N/Q/R/V/Y	5	1	0
AGG	Q	E/K/R	N/Q	H/K/R	R	D/E/K/N/Q/Y	5	1	0
CAA	E	D/K/N/Q/R/S/T/Y	N	D/E/H/K/Q/R/S/T/Y	Q	E/H/K/Q/R/Y	5	1	0
CAG	E	D/K/N/Q/R	Q	D/E/K/N/Q/R	R	D/E/H/K/N/Q/S/T/Y	5	1	0
CAT	E	D/K/N/Q/R/S/T/Y	N/Q	D/E/H/K/N/Q/R/S/T/Y	N/Q	D/E/H/K/Q/R/Y	5	1	0
CCC	E	D/K/N/Q/R	E/Q	D/H/K/N/S/T/Y	R	E/H/K/N/Q/S/T/Y	5	1	0
CCT	E	K/Q/R	D/N	E/H/K/N/Q/S/T	Q	E/H/K/N/Q/R/S/T/Y	5	1	0
GCA	R	H/K/N/Q/S/T/Y	D	E/H/K/N/Q/R/S/T/Y	Q	A/E/H/K/L/N/Q/R/S/T/V/Y	5	1	0
GCG	R	H/K/N/Q/S/T/Y	D/E	H/K/N/Q/S/T	R	D/E/H/K/N/Q/S/T/Y	5	1	0
GGA	R	H/K/N/Q/Y	N	H/K/Q/R/S/T/Y	Q	A/D/E/H/K/L/N/Q/R/V/Y	5	1	0
AAA	K/R		N	D/E/H/K/Q/R/Y	Q	D/E/H/K/L/N/Q/R/Y	5	0	0
AGC	Q	E/K/R	H	K/N/Q/R/S/T/Y	R	D/E/H/K/N/Q/Y	5	0	0
AGT	Q	E/K/R	H	K/N/Q/R/S/T	Q	D/E/H/K/L/N/Q/R/S/T/Y	5	0	0
GCG	R	H/K/N/Q/Y	H	K/N/Q/R/S/T/Y	R	D/E/H/K/N/Q/S/T/Y	5	0	0
GCG	R	H/K/N/Q/Y	H	K/N/Q/R/S/T/Y	R	D/E/H/K/N/Q/S/T/Y	5	0	0
GGT	R	H/K/N/Q/Y	H	K/N/Q/R/S/T	N/Q	D/E/K/N/Q/R/S/T	5	0	0
CAC	E	K/Q/R	N	D/E/K/Q/R	E	H/K/Q/R/Y	4	2	0
CCA	E	K/Q/R	D	E/K/N/Q/S/T	Q	A/E/H/K/L/N/Q/R/S/T/V/Y	4	2	0
CCG	E	D/K/N/Q/R	D	E/H/K/N/Q/R/S/T/Y	R	D/E/H/K/N/Q	4	2	0
CQA	E	D/K/N/Q/R	N	H/K/Q/R/S/T/Y	Q	A/D/E/H/K/L/N/Q/R/S/T/V/Y	4	2	0
CGC	E	D/K/N/Q/R/S/T/Y	H	K/N/Q/R/S/T/Y	R	D/E/H/K/N/Q/Y	4	1	0
CGG	E	D/K/N/Q/R	H	K/N/Q/R/Y	R	D/E/H/K/N/Q/S/T/Y	4	1	0
CGT	E	D/K/Q/R	H	K/N/Q/R/Y	N/Q	D/E/H/K/N/Q/R/Y	4	1	0

Table 3

ATA	Q	E/K/R	A/I/L/V	H/K/N/Q/R/S/T/Y	Q	E/H/I/K/L/Q/R/V/Y	4	0	1
ATC	Q	D/E/K/N/H	I/L	D/E/H/K/N/Q/R/S/T/V/Y	E/R	H/K/N/Q/R/Y	4	0	1
ATT	Q	E/H/K/R	I/L/V	D/E/H/K/N/Q/R/S/T/Y	R	D/E/H/K/N/Q/Y	4	0	1
GTA	R	E/K/R	I/L/V	H/K/N/Q/R/S/T/Y	Q	D/E/H/K/N/Q/R/S/T/Y	4	0	1
GTC	R	H/K/Q/Y	A/I/L/V	D/E/H/K/N/Q/R/S/T/Y	Q	A/D/E/H/I/K/L/N/Q/R/S/T/V/Y	4	0	1
GTG	R	H/K/Q	I/L	H/K/N/Q/R/S/T/V	E/R	H/K/N/Q/S/T/Y	4	0	1
GTT	R	H/K/N/Q/Y	I/L/V	H/K/N/Q/S/T	R	D/E/H/K/N/Q/S/T/Y	4	0	1
TAA	I/L	K/Q/R/V	N	A/H/K/N/Q/R/S/T/Y	Q	D/E/H/I/K/L/N/Q/R/V/Y	4	0	1
TAT	I/L/V	K/Q/R	N	D/E/H/K/Q/R/Y	Q	A/E/H/I/K/L/N/Q/R/V/Y	4	0	1
TTC	I/L	K/N/Q/R/Y	N	D/E/H/K/N/Q/R/S/T/Y	R	D/E/H/K/N/Q/Y	4	0	1
TCT	I/L	A/H/K/N/Q/R/V	E/Q	D/H/K/N/Q/R/S/T/Y	Q	D/E/H/I/K/L/N/Q/R/S/T/V/Y	4	0	1
		H/K/Q/R	D/N	E/H/K/N/Q/R/S/T/Y	N/Q	D/E/H/K/Q/R/Y	4	0	1
CTA	E	Q	A/I/L/V	H/K/N/Q/S/T	Q	A/E/I/K/L/N/Q/R/S/T/V	3	1	1
CTC	E	D/K/N/Q/R	I/L	H/K/N/Q/R/S/T/V/Y	E/R	D/H/K/N/Q/Y	3	1	1
CTG	E	D/K/N/Q/R	I/L/V	K/N/Q/S/T	R	D/E/K/N/Q/Y	3	1	1
CTT	E	D/K/N/Q/R	I/L	D/E/H/K/N/Q/R/S/T/V/Y	Q	D/E/H/K/N/Q/R/Y	3	1	1
TCA	I/L	K/Q/R	D	E/H/K/N/Q/R/S/T/Y	Q	A/E/H/I/K/L/N/Q/R/S/T/V/Y	3	1	1
TCG	I/L	K/Q/R	D	E/H/K/N/Q/R/S/T/Y	R	D/E/H/K/N/Q/S/T/Y	3	1	1
TGA	I/L	K/Q/R	N	H/K/Q/R/S/T/Y	Q	D/E/H/K/N/Q/R/Y	3	1	1
TAC	I/L/V	K/Q/R	N	D/E/H/K/Q/R/Y	E	H/K/N/Q/R/Y	3	0	1
TGC	I/L/V	H/K/N/Q/R/S/T	H	K/N/Q/R/S/T/Y	R	D/E/H/K/N/Q/Y	3	0	1
TGG	I/L	K/Q/R	H	K/N/Q/R/Y	R	D/E/H/K/N/Q/S/T/Y	3	0	1
TGT	I	H/K/L/N/Q/R	H	K/N/Q/R/S/T/Y	Q	E/H/I/K/L/N/Q/R/Y	3	0	1
TTA	I/L	K/N/Q/R	A/I/L/V	D/E/H/K/N/Q/R/S/T/Y	Q	A/E/H/I/K/L/N/R/S/T/V/Y	2	0	2
TTC	I/L/V	K/N/Q/R	I/L	A/H/K/N/Q/R/S/T/V/Y	E/R	D/H/K/N/Q/S/T/Y	2	0	2
TTG	I/L	K/Q/R/Y	I/L/V	H/K/N/Q/S/T	R	D/E/H/K/N/Q/Y	2	0	2
TTT	I/L	K/N/Q/R/V	I/L/V	A/H/K/N/Q/R/S/T/Y	Q	D/E/H/K/N/Q/R/Y	2	0	2

where / separates alternative amino acids

The results of the molecular modeling analysis of various ZF  $\alpha$ -helix complexes with the 64 different DNA triplets (Tables 1, 2 and 3), and the findings of spatial peculiarities in the three contact positions, are reflected in the ZF-DNA recognition rules. On the basis of the rules set forth in Tables 1, 2 and 3, DBP's with optimal binding affinity for any target DNA sequence can be designed. The "Column A" designations, i.e., the "A Rules," in Tables 1-3, show the amino acids with optimal binding for a given codon (triplet). The "Column B" designations, i.e., the "B rules," in Tables 1 and 3, show the amino acids with secondary, but still significant, binding affinity for a given triplet.

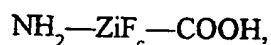
The column A rules range from the strongest triplet recognition with six H-bonds, zero water contracts and zero hydrophobic contacts with an evaluated energy of  $(5 \times 6) + (2 \times 0) + (1 \times 0) = 30$  to two hydrogen bonds, zero water contacts and two hydrophobic contacts with an evaluated energy of  $(5 \times 2) + (2 \times 0) + (1 \times 2) = 12$ . The Column A rules ordinarily have a choice of just one or two amino acids in positions  $Z_1$ ,  $Z_2$  and  $Z_3$ . The column B rules, by comparison, have from three possible amino acids in each of the  $Z_1$ ,  $Z_2$  and  $Z_3$  positions to as many as eighteen amino acids in different contacting arrangements in each of the  $Z_1$ ,  $Z_2$  and  $Z_3$  positions. In the evaluation of the column B energies, there are a large number different groupings of three amino acids in positions  $Z_1$ ,  $Z_2$  and  $Z_3$ . The minimum energy is three hydrogen bonds, zero water contacts and zero hydrophobic contacts with an evaluated energy of  $(5 \times 3) + (2 \times 0) + (1 \times 0) = 15$ . The maximum energy evaluation for these combinations is, on average, three hydrogen bonds and either two water contacts or two hydrophobic contacts, with an evaluated energy of from  $(5 \times 3) + (2 \times 2) + (1 \times 0) = 19$  down to  $(5 \times 3) + (2 \times 0) + (1 \times 2) = 17$ . Thus, the column B rules have a narrower energy range (i.e., from 19 down to 15) than do the column A rules, which have an energy range from 30 down to 12. The narrow energy range for the column B rules means that the 64 different rules do not distinguish on the basis of energy as well as the 64 column A rules.

For example, as set forth in Table 2, a DBP which binds optimally to the DNA base triplet guanine-cytosine-cytosine (GCC) is one wherein the portion of the protein responsible for the binding to the triplet is a ZF domain within which is contained a segment having the sequence  $Z_3XXZ_2LXZ_1H$ , wherein  $Z_1$  is an arginine which interacts with position 1 of the DNA triplet;  $Z_2$  is a glutamine or a glutamic acid which interacts with position 2 of the DNA triplet;  $Z_3$  is an arginine which interacts with position 3 of the DNA triplet; X is an arbitrary amino acid; L is leucine and H is histidine.

As set forth in Table 1 or 3 (see the "column B" entries for the  $Z_1$ ,  $Z_2$ , and  $Z_3$  positions for a given codon), a DBP which effectively, if not optimally, binds to the DNA base triplet guanine-cytosine-cytosine (GCC) is one wherein the portion of the protein responsible for the binding to the triplet is a ZF domain within which is contained a segment having the sequence  $Z_3XXZ_2LXZ_1H$ , wherein  $Z_1$  is an amino acid selected from the group consisting of histidine, lysine, glutamine, asparagine, tyrosine, serine and threonine which interacts with position 1 of the DNA triplet;  $Z_2$  is an amino acid selected from the group consisting of glutamine, asparagine, aspartic acid, serine, threonine, arginine, histidine, and lysine which interacts with position 2 of the DNA triplet;  $Z_3$  is an amino acid selected from the group consisting of glutamine, asparagine, glutamic acid, aspartic acid, histidine, lysine, tyrosine, serine and threonine which interacts with position 3 of the DNA triplet; X is an arbitrary amino acid; L is leucine and H is histidine.

It will be appreciated, of course, that DBP's of intermediate affinity, i.e., ones wherein the  $Z_1$ ,  $Z_2$  and  $Z_3$  contact amino acids are selected according to a combination of the "A" and "B Rules," can be designed. For example, in the segment  $Z_3XXZ_2LXZ_1H$  within a ZF domain for binding to the triplet GCC,  $Z_1$  could be an arginine;  $Z_2$  could be a glutamine or a glutamic acid; and  $Z_3$  could be selected from the group consisting of glutamine, asparagine, glutamic acid, aspartic acid, histidine, lysine, tyrosine, serine and threonine.

The basic building block for such proteins is denoted by the formula:



where  $ZiF_c$  is a ZF domain of the form

$Y/FXCX_{2-4}CG/DK/RXFXZ_3XXZ_2LXZ_1HX_{3-5}H$ , where

$Z_1$ ,  $Z_2$ , and  $Z_3$  are amino acids chosen from Table 1, 2 or 3 to correspond to the three bases of the DNA triplet, and the remaining components of the formula are as described earlier in the description of Formula I.

In the preferred embodiment of the invention, a zinc-finger domain for binding to a given DNA triplet is designed by selection of the appropriate AA's in Table 2 or in column A

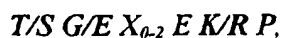
of Table 1 or Table 3. In another embodiment of the invention, the ZF domain is designed by selection from among the AA's set forth for a given DNA triplet in column B of Table 1 or 3.

One such domain is required for each triplet of the target sequence; for a target string of only 3 bases, the above formula defines the protein.

If the target string of DNA is 6 bases, the DBP design is extended as follows:

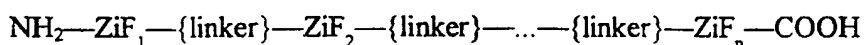


where  $\text{ZiF}_1$  and  $\text{ZiF}_2$  are ZF domains designed, as shown above for  $\text{ZiF}_e$ , to bind to the first and second triplets of the six bases, and {linker} is an amino acid sequence conforming to the pattern

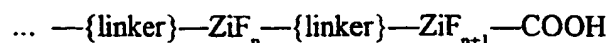
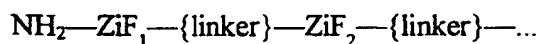


again wherein the components are as defined previously in Formula I.

If 1) the target string of DNA contains 9, 12, or a higher multiple of 3 bases; 2) it is required to design a DBP for  $3n+3$  bases; and 3) the DBP for the first  $3n$  bases is given by the sequence:



then the DBP design is extended recursively and the required DBP is specified by the sequence:



where  $\text{ZiF}_{n+1}$  is a ZF domain designed, as shown above for  $\text{ZiF}_e$ , to bind with the  $n^{\text{th}}+1$  triplet of the target sequence of base pairs.

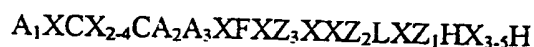
Figure 3 provides a schematic representation of a ZF protein wherein  $n=3$ , i.e., one which has 3 ZF domains (i.e.,  $n=3$ ) connected by linker sequences and is designed to bind to a target DNA string of 9 ( $3n$ ) bases.

The above rules enable ready determination of the optimal amino acid(s) for binding to any given DNA triplet and thus the identification and positioning of the 3 amino acids in a ZF domain which would be the ideal component of a DBP for binding to the DNA triplet.

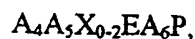
The application of the rules can then be extended to design of a DBP containing a set number,  $n_d$ , of ZF domains, which DBP binds to a target stretch of  $3n_d$  nucleotides within a given DNA sequence. The target  $3n_d$  stretch of nucleotides, and the collection and order of  $n_d$  domains in the DBP, are such that the binding energy for the DBP and target DNA sequence is the highest possible for any pairing of a DBP containing the set number,  $n_d$ , of ZF domains with any stretch of  $3n_d$  nucleotides within the entire DNA molecule being screened.

Accordingly, the embodiment of the invention of primary importance is a method for designing such a DBP for a DNA sequence of any length. The method employs the rules disclosed above in combination with a means of screening and ranking all possible segments of  $3n_d$  nucleotides within the sequence by their affinities for DBP's containing  $n_d$  ZF domains to determine a unique DBP with the desired properties.

More particularly, the invention is directed to a method for designing a DBP, with multiple ZF domains connected by linker sequences, that binds selectively to a target DNA sequence within a given gene, each of said ZF domains having the formula



and each of said linkers having the formula



wherein

- (i) X is any amino acid; (ii)  $X_{2-4}$  is a peptide from 2 to 4 amino acids in length; (iii)  $X_{3-5}$  is a peptide from 3 to 5 amino acids in length; (iv)  $X_{0-2}$  is a peptide from 0 to 2 amino acids in length; (v)  $A_1$  is selected from the group consisting of phenylalanine and tyrosine; (vi)  $A_2$  is selected from the group consisting of glycine and aspartic acid; (vii)  $A_3$  is selected from the group consisting of lysine and arginine; (viii)  $A_4$  is selected from the group consisting of



threonine and serine; (viii)  $A_5$  is selected from the group consisting glycine and glutamic acid; (ix)  $A_6$  is selected from the group consisting of lysine and arginine; (x) C is cysteine; (xi) F is phenylalanine; (xii) L is leucine; (xiii) H is histidine; (xiv) E is glutamic acid; (xv) P is proline; and (xvi)  $Z_1$ ,  $Z_2$  and  $Z_3$  are the base-contacting amino acids, comprising the steps of:

- (a) setting a genome to be screened;
- (b) selecting the target DNA sequence in the genome for binding;
- (c) setting the number of zinc-finger domains to  $n_d$ ;
- (d) dividing the target DNA sequence into nucleotide blocks wherein each block contains  $n_z$  nucleotides using a first routine where  $n_z$  is determined using the following relationship:

$$n_z = 3n_d;$$

(e) assigning base-contacting amino acids at  $Z_1$ ,  $Z_2$  and  $Z_3$  to each ZF domain, according to the A Rules and /or B Rules set forth in Tables 1-3, of a DBP which binds to the first nucleotide block from step (d) as numbered from the first 5' nucleotide of the target gene sequence to generate a block-specific DBP and calculating the binding energy, Binding Energy<sub>block</sub>, of each ZF domain of each such block-specific DBP as the product of the binding energies, Binding Energy<sub>domain</sub>, of all zinc-finger domains of the polypeptide, each determined using the formula:

$$\text{Binding Energy}_{\text{domain}} = (5 \times \text{the number of hydrogen bonds}) + (2 \times \text{the number of H}_2\text{O contacts}) + (\text{the number of hydrophobic contacts});$$

- (f) subdividing the DBP from step (d) into blocks using a second routine to generate a subdivided DBP having three ZF domains;
- (g) screening the subdivided DBP from step (f) against the genome using a third routine to determine the number of binding sites in the genome for each subdivided DBP in the genome and assigning a binding energy for each such site using the following formula:

$$\text{Binding Energy}_{\text{site } n} = (5 \times \text{the number of hydrogen bonds}) + (2 \times \text{the number of H}_2\text{O contacts}) + (\text{the number of hydrophobic contacts});$$

- (h) calculating a ratio of binding energy,  $R_b$ , using a fourth routine for each nucleotide-block-specific DBP from step (e) using the following formula:

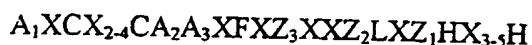
$R_b$  = Binding Energy<sub>block</sub> / the sum of all Binding Energy<sub>site n</sub>'s for all subdivided DBP's from step (g);

- (i) repeating steps (f) through (h) for each subdivided DBP wherein  $n_d \geq 4$ ;
- (j) repeating steps (d) through (i) for each nucleotide block in the target DNA sequence containing  $n_z$  nucleotides;
- (k) rank-ordering  $R_b$  numerical values obtained from step (h); and
- (l) selecting a DBP with an acceptable  $R_b$  value.

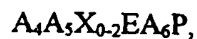
Preferred embodiments of this aspect of the invention are:

- 1) the design method as set forth above wherein the DBP  $R_b$  numerical value is the highest numerical value for all DBP's in step (h) that bind to the target DNA sequence.
- 2) the method above wherein the DBP  $R_b$  numerical value determined in step (h) is at least 10,000.
- 3) the method above wherein the number of ZF domains,  $n_d$ , is nine.
- 4) the method above wherein the rules for assigning base-contacting amino acids at  $Z_1$ ,  $Z_2$  and  $Z_3$  for each nucleotide block in step (e) are selected from rule set A.

The invention is further directed to a computer system for designing a DBP, with multiple ZF domains connected by linker sequences, that binds selectively to a target DNA sequence within a given gene, each of said ZF domains having the formula



and each of said linkers having the formula



wherein

- (i) X is any amino acid; (ii)  $X_{2-4}$  is a peptide from 2 to 4 amino acids in length; (iii)  $X_{3-5}$  is a peptide from 3 to 5 amino acids in length; (iv)  $X_{0-2}$  is a peptide from 0 to 2 amino acids in length; (v)  $A_1$  is selected from the group consisting of phenylalanine and tyrosine; (vi)  $A_2$  is selected from the group consisting of glycine and aspartic acid; (vii)  $A_3$  is selected from the group consisting of lysine and arginine; (viii)  $A_4$  is selected from the group consisting of threonine and serine; (ix)  $A_5$  is selected from the group consisting glycine and glutamic

acid; (ix) A<sub>6</sub> is selected from the group consisting of lysine and arginine; (x) C is cysteine; (xi) F is phenylalanine; (xii) L is leucine; (xiii) H is histidine; (xiv) E is glutamic acid; (xv) P is proline; and (xvi) Z<sub>1</sub>, Z<sub>2</sub> and Z<sub>3</sub> are the base-contacting amino acids, comprising the steps of:

- (a) setting a genome to be screened;
- (b) selecting the target DNA sequence in the genome for binding;
- (c) setting the number of ZF finger domains to n<sub>d</sub>;
- (d) dividing the target DNA sequence into nucleotide blocks wherein each block contains n<sub>z</sub> nucleotides using a first routine where n<sub>z</sub> is determined using the following relationship:

$$n_z = 3n_d;$$

- (e) assigning base-contacting amino acids at Z<sub>1</sub>, Z<sub>2</sub> and Z<sub>3</sub> to each ZF domain, according to the A Rules and/or B Rules set forth in Tables 1-3, of a DBP which binds to the first nucleotide block from step (d) as numbered from the first 5' nucleotide of the target gene sequence to generate a block-specific DBP and calculating the binding energy, Binding Energy<sub>block</sub>, of each ZF domain of each such block-specific DBP as the product of the binding energies, Binding Energy<sub>domain</sub>, of all domains of the DBP, each determined using the formula:

$$\text{Binding Energy}_{\text{domain}} = (5 \times \text{the number of hydrogen bonds}) + (2 \times \text{the number of H}_2\text{O contacts}) + (\text{the number of hydrophobic contacts});$$

- (f) subdividing the DBP from step (d) into blocks using a second routine to generate a subdivided DBP having three ZF domains;
- (g) screening the subdivided DBP from step (f) against the genome using a third routine to determine the number of binding sites in the genome for each subdivided DBP in the genome and assigning a binding energy for each such site using the following formula:

$$\text{Binding Energy}_{\text{site } n} = (5 \times \text{the number of hydrogen bonds}) + (2 \times \text{the number of H}_2\text{O contacts}) + (\text{the number of hydrophobic contacts});$$

- (h) calculating a ratio of binding energy, R<sub>b</sub>, using a fourth routine for each nucleotide block-specific DBP from step (e) using the following formula:

$$R_b = \text{Binding Energy}_{\text{block}} / \text{the sum of all Binding Energy}_{\text{site } n} \text{'s for all subdivided DBP's from step (g);}$$

- (i) repeating steps (f) through (h) for each subdivided DBP wherein  $n_d \geq 4$ ;
- (j) repeating steps (d) through (i) for each nucleotide block in the target DNA sequence containing  $n_z$  nucleotides;
- (k) rank-ordering  $R_b$  numerical values obtained from step (h);
- (l) selecting a DBP with an acceptable  $R_b$  value.

According to the instant invention,  $R_b$ , as defined in (h) above for both the design method and computer system, has a lower limit of 10,000. Preferably  $R_b$  is greater than  $10^6$ .

Preferred embodiments of this aspect of the invention are:

- 1) the computer system as set forth above wherein the DBP  $R_b$  numerical value is the highest numerical value for all DBP's in step (h) that bind to the target DNA sequence.
- 2) the computer system above wherein the DBP  $R_b$  numerical value determined in step (h) is at least 10,000.
- 3) the computer system above wherein the number of ZF domains,  $n_d$ , is nine.
- 4) the computer system above wherein the rules for assigning base-contacting amino acids at  $Z_1$ ,  $Z_2$  and  $Z_3$  for each nucleotide block in step (e) are selected from rule set A.

The method and computer system of the instant invention are further illustrated by the block flow diagrams of Figures 4-9.

Figure 4 shows the components of the computer system on which the DBP design process is implemented. A Central Processor Digital Computer (1) of any manufacture is provided with a Computer Program (2) written by the inventors. This Computer Program (2) reads a series of files described as DNA-Triple Energy Rules (6), Genome Descriptors (9), Genomic DNA Sequence (10) and Gene Features (5). The Central Processor (1) transforms this information into the DBP Blocking Fragment Files (7) and the Optimal DBP Designs for Genome (8).

Figure 5 shows that the Computer Program (2) in Figure 4 has two portions. The genomic data is first transformed by the Process Genome into Blocking Fragment Files function (2). These files are then used by the Design DBP's for a Genome function (3).

The Process Genome into Blocking Fragment Files block (2) of Figure 5 is represented in greater detail in Figure 6. For every  $n_d$  from 11 down to 3 the Genome Descriptors file (12) and the Genome DNA Sequence file (32) are read and transformed into the Unsorted Fragment File (7). This same Unsorted Fragment File (14) is transformed by the Sort function

(13) provided by the computer manufacturer into the Sorted Fragment file (15). The same Sorted Fragment File (30) is read and transformed eventually into the DBP-Size Blocking File (22).

The Design DBP's for a Genome block (3) of Figure 5 is represented in greater detail in Figure 7. The Genome Descriptors file (3), the Gene Features file (7), the Genome DNA Sequence file (9) and the DBP-Size Blocking Files (37) corresponding to the  $n_d$ 's from 11 down to 3 are read and used to transform the genomic DNA first into genes and then into a file of the Optimal DBP Designs for a Genome (38). The transformation and design process is done for all the genes in a genome.

The "Determine if Current-Sub-Window is in Current-Blocking-File" block (22) in Figure 7 is expanded in greater detail in Figure 8.

The "Calculate Binding-Energy-of-Blocking-Fragment" block (24) in Figure 7 is expanded in greater detail in Figure 9.

By applying the algorithm to a variety of DBP's of varying  $n_d$ , it was experimentally determined that a value for  $n_d$  of 9 is the best starting point in the algorithm, i.e., the process should begin with the search for 9-finger DBP's. This can be better understood in terms of the selection criterion,  $R_b$ , used in evaluating various DBP's. In short DBP's, e.g., ones wherein  $n_d = 4$  or 5, Binding Energy<sub>block</sub>, which increases geometrically as the product of all Binding Energy<sub>domain</sub>'s, is significantly lower, and Binding Energy<sub>site n</sub> values are relatively large. However, as  $n_d$  increases, the numerator of  $R_b$  increases dramatically, while, it has been observed, the denominator, representing "background" or "noise," does not significantly change. Thus, the case of  $n_d = 9$  provides assurance of high affinity and specificity of binding without also bringing on the possibility of undue computational needs.

However, it should also be emphasized that the present invention is not limited to the design of DBP's wherein  $n_d \leq 9$ . For that matter, it will also be appreciated that, while  $n_d = 9$  has been found to be the best starting point, the best DBP for a given situation may turn out to be one wherein  $n_d < 9$ , the length of the target DNA sequence notwithstanding. The concept of the invention can be applied to the design of DBP's of any length as required.

In any event, for a given DNA sequence of  $N$  nucleotides, there are  $N - 27$ , 9-finger DBP sequences. Each of these can be ordered in terms of strength of binding by evaluating the energy function for each 3-nucleotide segment as set forth in part (e) of the design method disclosed above.

In initial computational experiments, a selectable sequence could have no 8, 7, 6, 5, and 4-finger subsites; however, with the present system, only the sum of the subsite binding energies need be minimized. As a result, it does not matter whether the subsite binding energy comes from 3-finger subsites, 4-finger subsites or even (in principle) larger subsites. This simple change from logical exclusion to energetic exclusion has been mandated not so much by examination of the yeast genome, but more by examination of the worm genome.

The central portion of the instant algorithm is, in the case of finding an acceptable  $n_4$ -finger site (e.g., a 27-base segment for a 9-finger DBP), the search against all other  $n_4$ -finger sites in the entire genome to see if there are any similar sites. If such turns out to be the case, the DBP with the highest  $R_b$  value is selected. Furthermore, the algorithm checks to see if there are any equivalent 8-finger, 7-finger, 6-finger, 5-finger and 4-finger subsites in the whole genome for a given 9-finger site. In the event no acceptable 9-finger site is found, the algorithm then searches for a suitable 8-finger site. If necessary, the search is continued for a 7-finger site and so on, until an acceptable DBP binding site is found.

Within the search for a 9-finger DBP, the algorithm looks at all 27-base sequences, which are called "frames." Each frame is evaluated to determine its interaction with DNA and the interaction of all other subframes down to 3-finger subsites. The number of instances of each frame and subframe in the genome has been recorded during the genome processing phase of the execution of the software. The sequence of the frame or subframe is evaluated as a product of the binding energy of each ZF. Each ZF domain recognizes three DNA bases. The underlying DNA sequence that a ZF recognizes determines how many hydrogen bonds, water contacts and hydrophobic contact exist between the ZF and the DNA.

The way the algorithm detects whether a given  $n_2$ -base site occurs in other places in the genome is by looking in a B-tree for the site. The whole genome is processed for each of the  $n_4$ -finger sites. The algorithm contains means for sorting and merging the myriad fragments and, in the end, there is produced an ordered list of all the blocking fragments for all the different finger sizes.

### Example 1

The following is given as an example of how the search for, and design of, a DBP is typically carried out. It involves screening for 9-finger DBP's (i.e.,  $n_4 = 9$ ) to bind to a target DNA sequence of 100 nucleotides (i.e.,  $N = 100$ ). The sequence is screened, beginning with

position 1, for every 27-nucleotide sequence, i.e., 1-27, 2-28, 3-29 etc., in the entire 100-nucleotide sequence. Once this has been done, the 9-fingers are broken down into 3-finger sections, i.e., 1-3, 2-4, 3-5 etc. The algorithm scans and looks for relative strengths of binding. The idea is to maximize the ratio of DBP binding to subsite binding,  $R_b$ , thus eliminating those 9-mers interacting with the greatest numbers of subsites.

The algorithm of the present invention was applied to the genomes of *S. cerevisiae* and *C. elegans* as illustrated by the following examples:

#### Example 2

The algorithm has been applied to the screening of the yeast genome. Two chromosomes of yeast, containing 110 and 447 genes, respectively, have been processed. For each gene the algorithm selected the  $n_d$ -finger sequence with the lowest sum of subsite binding energies. In yeast the number of 3-finger blocking fragments is almost maximal (i.e.,  $4^9$ , versus  $4^{12}$  maximal). In the worm genome (see Example 3), the 3-finger blocking sequences are absolutely maximal. In yeast the 4-finger blocking sequences are large in number but the population of 5-finger blocking sequences is relatively small. In worm the 4-finger blocking sequences are larger in number than the 5-finger blocking sequences, but the latter are larger in number relative to yeast. In going in the future from worm to human, one can expect that the 4-finger blocking sequences might come close to saturation (i.e. close to  $4^{12}$ ).

The algorithmic analysis was performed for 2 of the 16 chromosomes of yeast. The 557 genes in the first two chromosomes seem to present a realistic picture of properties of all the chromosomes in the yeast genome. Sample calculations have been run on the whole yeast genome but these results are not different from those produced by calculating the properties of just two chromosomes' worth of genes. The results of the analysis of 100 yeast genes, typical of the findings throughout the analysis of the yeast genome, are presented in Table 4.

The power of the algorithm is further demonstrated in the results displayed in Figures 10-14. The figures display results obtained for all 557 genes of the two yeast chromosomes on which the studies were focused.

The strength of each acceptable 9-finger DBP can be calculated. Figure 10 shows that the strengths of binding of all the acceptable 9-finger DBP's are uniformly distributed. If this curve were bowed down, then the stronger frames would be more preferred. If this curve were bowed up, then the weaker frames would be preferred.

Figure 11 shows that the binding energies (Binding Energy<sub>block</sub>'s) of the acceptable 9-finger DBP's are uniformly distributed between  $10^{11}$  and  $10^{13}$  binding units.

Figure 12 shows that the distribution of the sum of the spurious subsite binding energies (Binding Energy<sub>site</sub>'s) is itself uniform in the range of  $10^6$  to  $10^8$  binding units.

Figure 13 is a nonlogarithmic version of Figure 12. It shows that most of the acceptable 9-finger DBP's have spurious subsite binding energies of less than  $5 \times 10^6$ .

Figure 14, produced by taking the ratios of the Figure 11 values to those of Figure 12, is a graph of the  $R_b$ 's for the 9-finger DBP's. This chart shows that the ratio of the DBP binding strength of the acceptable 9-finger DBP's to the sum of the binding energies of the spurious subsite interactions varies from  $10^4$  to  $10^6$ .

The analytical tools of the present invention were also employed in the further analysis of a single yeast gene, YAR073, in particular the 300-bp region of the promoter immediately upstream of the coding region. The full sums of the subsite binding energies (SBE's) for each 27-base frame in this portion of the gene were determined; the results are depicted graphically in Figure 15. The primary binding energies (BE's) were also determined, and a correlation was found between the SBE values and the values of the ratios of BE:SBE ( $R_b$ ). Still further (Figure 16), it was seen that the peaks of the plot of the  $R_b$  values correspond to the footprints of the transcription factors of the same gene (determined in a separate study).

### Example 3

Application of the algorithm according to the instant invention to 100 genes in *C. elegans* showed that the system can be applied as successfully to *C. elegans* as to *S. cerevisiae*. The results of analysis of the 100 *C. elegans* genes are presented in Table 5.

In Figure 17, it can be seen that, for one of the analyzed *C. elegans* genes, only a 5-finger DBP could be designed. For another gene, only a 7-finger DBP could be designed. These two genes, 2 and 32, are not seen in Table 5, since it presents results of the analysis only for those genes (98 out of 100) for which a 9-mer could be designed. In any event, the results depicted in Figure 17 are in keeping with the expectation for analysis of the entire *C. elegans* genome namely, that the distribution of 5- through 9-finger DBP's is somewhat different than in *S. cerevisiae*.



Figure 18 represents the same analysis for the *C. elegans* genes as was depicted in Figure 14 for *S. cerevisiae* genes. Figure 18 shows a similar  $R_b$  value distribution to that seen in Figure 14.

Examples 2 and 3 demonstrate the applicability of the instant invention to the design of DBP's for the genomes of two widely disparate organisms. The various results of the application of the algorithm to the yeast genome, in particular, and also to the worm genome, show the power of the algorithmic tool and demonstrate its foundation in reality, i.e., that it does not merely provide a random and/or theoretical analysis. It is to be expected, on the basis of these analyses, that the inventive algorithm can be extended to the design of DBP's for any desired segment of the genome of any organism of interest, including that of a human.

Although the instant algorithm involves a search against the entire genome of an organism, the results of the present studies strongly indicate that lack of complete knowledge of the genome of a given organism would not constitute an impediment to application of the present invention to the design of DBP's for that organism. One would expect to be able to use the knowledge of block sequences obtained in the studies presented herein on *S. cerevisiae* (a unicellular organism) and *C. elegans* (a multicellular organism) to form valid estimates of allowable sequences for the systems of higher eukaryotes.

For example, the present studies on yeast and worm indicate that the genomic "noise," in this context the spurious binding site energies, is relatively constant, and this can be projected to higher, more complex organisms as well. In other words, one would expect from the demonstrated combinatorics of DNA sequences to be able to extrapolate, or extend, the present algorithm to the analysis of more complex genomes, however much is known of the specific sequences therein, with the object of designing effective DBP's. Furthermore, as the entire genomes of larger organisms, e.g., *D. melanogaster*, become known, they will provide further keys to the analysis of the genomes of higher organisms, including humans.

A DBP as specified above may be built by using standard protein synthesis techniques; or, employing the standard genetic code, may be used as the basis for specifying and constructing a gene whose expression is the DBP.

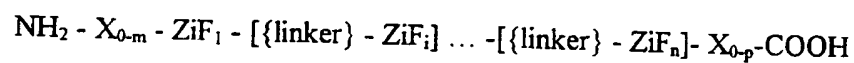
Proteins so designed can be used in any application requiring accurate and tight binding to a DNA target sequence. For example, a DBP, according to the instant invention, can be coupled with a DNA endonuclease activity. When the resultant molecule binds to the target DNA, said DNA can be cut at a fixed displacement from the DBP binding site.

Similarly, in instances in which the target DNA sequence is a promoter, one can produce a promoter-specific DBP which, when bound, will act to alter (i.e., enhance, attenuate or even terminate) expression of a given gene or, alternatively, genes under control of that promoter.

As another application, a DBP could be designed to bind specific DNA sequences when attached to solid supports. Such solid supports could include styrene beads, acrylamide well-plates or glass substrates.

In order to realize the specific applications mentioned above, as well as the full scope of applications possible through the instant invention, the DBP can be designed as set forth above to include the added feature of a pre- and/or postdomain amino acid sequence of arbitrary length. This would include, for example, the coupling of the basic DBP to an endonuclease or to a reporter or to a sequence by which the DPB could be coupled to a solid support.

Accordingly, the instant invention includes DBP's that bind to a predetermined target double-stranded DNA sequence of  $3n$  (where  $n \geq 1$ ) base pairs in length of the form:



wherein each  $\text{ZiF}_1$  to  $\text{ZiF}_n$  is a ZF domain of the form set forth above; {linker} is an amino acid sequence as set forth above;  $\text{X}_{0-m}$  stands for a sequence of from 0 to  $m$  amino acids and  $\text{X}_{0-p}$  stands for a sequence of from 0 to  $p$  amino acids. The values for  $m$  and  $p$  and the identities of the amino acids are determined by the particular protein(s) or amino acid sequence(s) to be coupled to the DBP for a given application.

In a further embodiment of the invention, the  $\text{Zn}^{+2}$  atom, which forms a complex with the two cysteine and two histidine amino acids in a specific ZF motif, can be substituted by a  $\text{Co}^{+2}$  or a  $\text{Cd}^{+2}$  atom, thus making a "cobalt finger" or a "cadmium finger."

The rules presented in Table 2 ("rule set A") are to be regarded as the "first choice" rules for optimal combinations in ZF-DNA recognition. However, it should be emphasized, as indicated in column B ("rule set B") of Table 1 or Table 3, that there are many alternative AA combinations that would also be expected to be important in the design of DNA-binding proteins capable of forming useful ZF-DNA complexes.

... a translation scheme for conflicts - map's designed against the coding region of each gene

[illegible]

TABLE 4

Scanning the complete *S. cerevisiae* genome for conflicts - DAP's designed against the coding region of each gene

[illegible]

**TABLE 4**  
**(cont'd.).**



[illegible]

**TABLE 5**  
**(cont'd.)**

## REFERENCES

1. Miller, J., McLachlan, A.D. and Klug, A. (1985) EMBO J. 4, 1609-1614.
2. Berg, J.M. (1988) Proc. Natl. Acad. Sci. USA 85, 99-102.
3. Gibson, T.J., Postma, J.P.M., Brown, R.S. and Argos, P. (1988) Protein Eng. 2, 209-218.
4. Lee, M.S., Gippert, G.P., Soman, K.V., Case, D.A. and Wright, P.E. (1989) Science 245, 635-637.
5. Nardelli, J., Gibson, T.J., Vesque, C. and Charnay, P. (1991) Nature 349, 175-178.
6. Thiesen, H.J. and Bach, C. (1991) FEBS Lett. 283, 23-26.
7. Pavletich, N.P. and Pabo, C.O. (1991) Science 252, 809-817.
8. Berg, J.M. (1992) Proc. Natl. Acad. Sci. USA 89, 11109-11110.
9. Kriwacki, R.W., Schultz, S.C., Steitz, T.A. and Caradonna, J.P. (1992) Proc. Natl. Acad. Sci. USA 89, 9759-9763.
10. Jacobs, G.H. (1992) EMBO J. 11, 4507-4517.
11. Gogos, J.A., Hsu, T., Bolton, J. and Kafatos, F.C. (1992) Science 257, 1951-1955.
12. Fairall, L., Harrison, S.D., Travers, A.A. and Rhodes, D. (1992) J. Mol. Biol. 226, 349-366.
13. Klevit, R.E. (1991) Science 253, 1367 and 1393.
14. Desjarlais, J.R. and Berg, J.M. (1992) Proc. Natl. Acad. Sci. USA 89, 7345-7349.
15. Rebar, E.J., and Pabo, C.O., U.S. Patent No. 5,789,538, August 4, 1998.
16. Beerli, R.R., Segal, D.J., Dreier, B., and Barbas, C.F. (1998) Proc. Natl. Acad. Sci. U.S.A. 95, 14628-14633.
17. Jordan, S.R. and Pabo, C.O. (1988) Science 242, 893-899.
18. Aggarwal, A.K., Rodgers, D.W., Drott, M., Ptashne, M. and Harrison S.C. (1988) Science 242, 899-907.
19. Letovsky, J., Dyran, W.S. (1989) Nucleic Acids Res. 17, 2639-2653.
20. Kissinger, C.R., Liu, B., Martin-Blanco, E., Kornberg, T.B. and Pabo, C.O. (1990) Cell 63, 579-590.
21. Wolberger, C., Vershon, A.K., Liu, B., Johnson, A.D. and Pabo, C.O. (1991) Cell 67, 517-528.
22. Seeman, N.C., Rosenberg, J.M. and Rich, A. (1976) Proc. Natl. Acad. Sci USA 73, 804-808.

23. Mikelsaar, R.-H., Bruskov, V.I. and Poltev, V.I. (1985) New precision space-filling atomic-molecular models, Pushchino.
24. Mikelsaar, R. (1986) Trends in Biotechnology 4, 162-163.

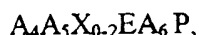


What is claimed is:

1. A method for designing a DBP, with multiple ZF domains connected by linker sequences, that binds selectively to a target DNA sequence within a given gene, each of said ZF domains having the formula



and each of said linkers having the formula



wherein

(i) X is any amino acid; (ii)  $X_{2-4}$  is a peptide from 2 to 4 amino acids in length; (iii)  $X_{3-5}$  is a peptide from 3 to 5 amino acids in length; (iv)  $X_{0-2}$  is a peptide from 0 to 2 amino acids in length; (v)  $A_1$  is selected from the group consisting of phenylalanine and tyrosine; (vi)  $A_2$  is selected from the group consisting of glycine and aspartic acid; (vii)  $A_3$  is selected from the group consisting of lysine and arginine; (viii)  $A_4$  is selected from the group consisting of threonine and serine; (ix)  $A_5$  is selected from the group consisting glycine and glutamic acid; (x)  $A_6$  is selected from the group consisting of lysine and arginine; (xi) C is cysteine; (xii) F is phenylalanine; (xiii) L is leucine; (xiv) H is histidine; (xv) E is glutamic acid; (xvi) P is proline; and (xvii)  $Z_1$ ,  $Z_2$  and  $Z_3$  are the base-contacting amino acids, which method comprises an algorithm comprising the steps of:

- (a) setting a genome to be screened;
- (b) selecting the target DNA sequence in the genome for binding;
- (c) setting the number of ZF domains to  $n_d$ ;
- (d) dividing the target DNA sequence into nucleotide blocks wherein each block contains  $n_z$  nucleotides using a first routine where  $n_z$  is determined using the following relationship:

$$n_z = 3n_d;$$

- (e) assigning base-contacting amino acids at  $Z_1$ ,  $Z_2$  and  $Z_3$  to each ZF domain, according to the A Rules and/or B Rules set forth in Tables 1-3 of the specification, of a DBP

which binds to the first nucleotide block from step (d) as numbered from the first 5' nucleotide of the target gene sequence to generate a block-specific DBP and calculating the binding energy,  $\text{Binding Energy}_{\text{block}}$ , of each ZF domain of each such block-specific DBP as the product of the binding energies,  $\text{Binding Energy}_{\text{domain}}$ , of all ZF domains of the DBP, each determined using the formula:

$$\text{Binding Energy}_{\text{domain}} = (5 \times \text{the number of hydrogen bonds}) + (2 \times \text{the number of H}_2\text{O contacts}) + (\text{the number of hydrophobic contacts});$$

- (f) subdividing the DBP from step (d) into blocks using a second routine to generate a subdivided DBP having three ZF domains;
- (g) screening the subdivided DBP from step (f) against the genome using a third routine to determine the number of binding sites in the genome for each subdivided DBP in the genome and assigning a binding energy for each such site using the following formula:

$$\text{Binding Energy}_{\text{site } n} = (5 \times \text{the number of hydrogen bonds}) + (2 \times \text{the number of H}_2\text{O contacts}) + (\text{the number of hydrophobic contacts});$$

- (h) calculating a ratio of binding energy,  $R_b$ , using a fourth routine for each nucleotide block-specific DBP from step (e) using the following formula:

$$R_b = \text{Binding Energy}_{\text{block}} / \text{the sum of all Binding Energy}_{\text{site } n} \text{'s for all subdivided DBP's from step (g);}$$

- (i) repeating steps (f) through (h) for each subdivided DBP wherein  $n_d \geq 4$ ;
  - (j) repeating steps (d) through (i) for each nucleotide block in the target DNA sequence containing  $n_z$  nucleotides;
  - (k) rank-ordering  $R_b$  numerical values obtained from step (h); and
  - (l) selecting a DBP with an acceptable  $R_b$  value.
2. The method of claim 1 wherein the DBP selected is that whose  $R_b$  numerical value is the highest numerical value for all DBP's in step (h) that bind to the target DNA sequence.
  3. The method of claim 1 wherein the DBP  $R_b$  numerical value determined in step (h) is at least 10,000.

4. The method of claim 1 wherein the number of ZF domains,  $n_d$ , is nine.
5. The method of claim 1 wherein the rules for assigning base-contacting amino acids at  $Z_1$ ,  $Z_2$  and  $Z_3$  for each nucleotide block in step (e) are selected from rule set A.
6. The method of claim 1 wherein the rules for assigning base-contacting amino acids at  $Z_1$ ,  $Z_2$  and  $Z_3$  for each nucleotide block in step (e) are selected from rule set B.
7. The method of claim 1 wherein rules for assigning base-contacting amino acids at  $Z_1$ ,  $Z_2$  and  $Z_3$  for each nucleotide block in step (e) are a combination selected from rule sets A and B.
8. A computer system for designing a DBP, with multiple ZF domains connected by linker sequences, that binds selectively to a target DNA sequence within a given gene, each of said ZF domains having the formula



and each of said linkers having the formula



wherein

(i) X is any amino acid; (ii)  $X_{2-4}$  is a peptide from 2 to 4 amino acids in length; (iii)  $X_{3-5}$  is a peptide from 3 to 5 amino acids in length; (iv)  $X_{0-2}$  is a peptide from 0 to 2 amino acids in length; (v)  $A_1$  is selected from the group consisting of phenylalanine and tyrosine; (vi)  $A_2$  is selected from the group consisting of glycine and aspartic acid; (vii)  $A_3$  is selected from the group consisting of lysine and arginine; (viii)  $A_4$  is selected from the group consisting of threonine and serine; (ix)  $A_5$  is selected from the group consisting of glycine and glutamic

acid; (ix) A<sub>6</sub> is selected from the group consisting of lysine and arginine; (x) C is cysteine; (xi) F is phenylalanine; (xii) L is leucine; (xiii) H is histidine; (xiv) E is glutamic acid; (xv) P is proline; and (xvi) Z<sub>1</sub>, Z<sub>2</sub> and Z<sub>3</sub> are the base-contacting amino acids, which computer system comprises means for design which include an algorithm comprising the steps of:

- (a) setting a genome to be screened;
- (b) selecting the target DNA sequence in the genome for binding;
- (c) setting the number of ZF domains to n<sub>d</sub>;
- (d) dividing the target DNA sequence into nucleotide blocks wherein each block contains n<sub>z</sub> nucleotides using a first routine where n<sub>z</sub> is determined using the following relationship:

$$n_z = 3n_d;$$

(e) assigning base-contacting amino acids at Z<sub>1</sub>, Z<sub>2</sub> and Z<sub>3</sub> to each ZF domain, according to the A Rules and/or B Rules set forth in Tables 1-3 of the specification, of a DBP which binds to the first nucleotide block from step (d) as numbered from the first 5' nucleotide of the target gene sequence to generate a block-specific DBP and calculating the binding energy, Binding Energy<sub>block</sub>, of each ZF domain of each such block-specific DBP as the product of the binding energies, Binding Energy<sub>domain</sub>, of all ZF domains of the DBP, using the formula:

$$\text{Binding Energy}_{\text{domain}} = (5 \times \text{the number of hydrogen bonds}) + (2 \times \text{the number of H}_2\text{O contacts}) + (\text{the number of hydrophobic contacts});$$

- (f) subdividing the DBP from step (d) into blocks using a second routine to generate a subdivided DBP having three ZF domains;
- (g) screening the subdivided DBP from step (f) against the genome using a third routine to determine the number of binding sites in the genome for each subdivided DBP in the genome and assigning a binding energy for each such site using the following formula:

$$\text{Binding Energy}_{\text{site } n} = (5 \times \text{the number of hydrogen bonds}) + (2 \times \text{the number of H}_2\text{O contacts}) + (\text{the number of hydrophobic contacts});$$

- (h) calculating a ratio of binding energy, R<sub>b</sub>, using a fourth routine for each nucleotide block-specific DBP from step (e) using the following formula:

$$R_b = \text{Binding Energy}_{\text{block}} / \text{the sum of all Binding Energy}_{\text{site } n} \text{'s for all subdivided DBP's from step (g);}$$

- (i) repeating steps (f) through (h) for each subdivided DBP wherein  $n_d \geq 4$  ;
- (j) repeating steps (d) through (i) for each nucleotide block in the target DNA sequence containing  $n_z$  nucleotides;
- (k) rank-ordering  $R_b$  numerical values obtained from step (h); and
- (l) selecting a DBP with an acceptable  $R_b$  value.

9. The computer system according to claim 8 wherein the DBP selected is that whose  $R_b$  numerical value is the highest numerical value for all DBP's in step (h) that bind to the target DNA sequence.

10. The computer system according to claim 8 wherein the DBP  $R_b$  numerical value determined in step (h) is at least 10,000.

11. The computer system according to claim 8 wherein the number of ZF domains,  $n_d$ , is nine.

12. The computer system according to claim 8 wherein the rules for assigning base-contacting amino acids at  $Z_1$ ,  $Z_2$  and  $Z_3$  for each nucleotide block in step (e) are selected from rule set A.

13. The computer system according to claim 8 wherein the rules for assigning base-contacting amino acids at  $Z_1$ ,  $Z_2$  and  $Z_3$  for each nucleotide block in step (e) are selected from rule set B.

14. The computer system according to claim 8 wherein the rules for assigning base-contacting amino acids at  $Z_1$ ,  $Z_2$  and  $Z_3$  for each nucleotide block in step (e) are a combination selected from rule sets A and B.



FIGURE 2

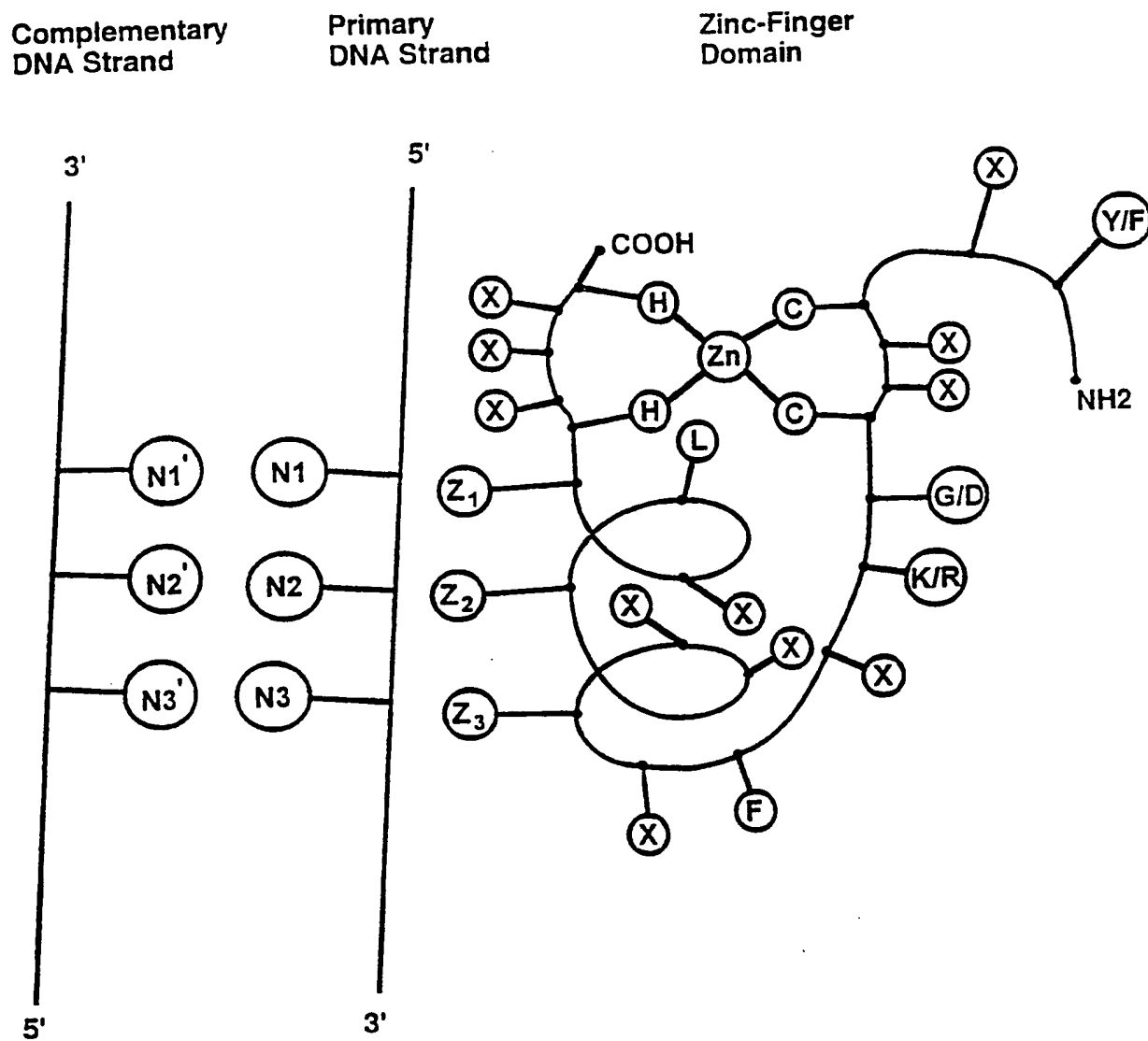






FIGURE 4

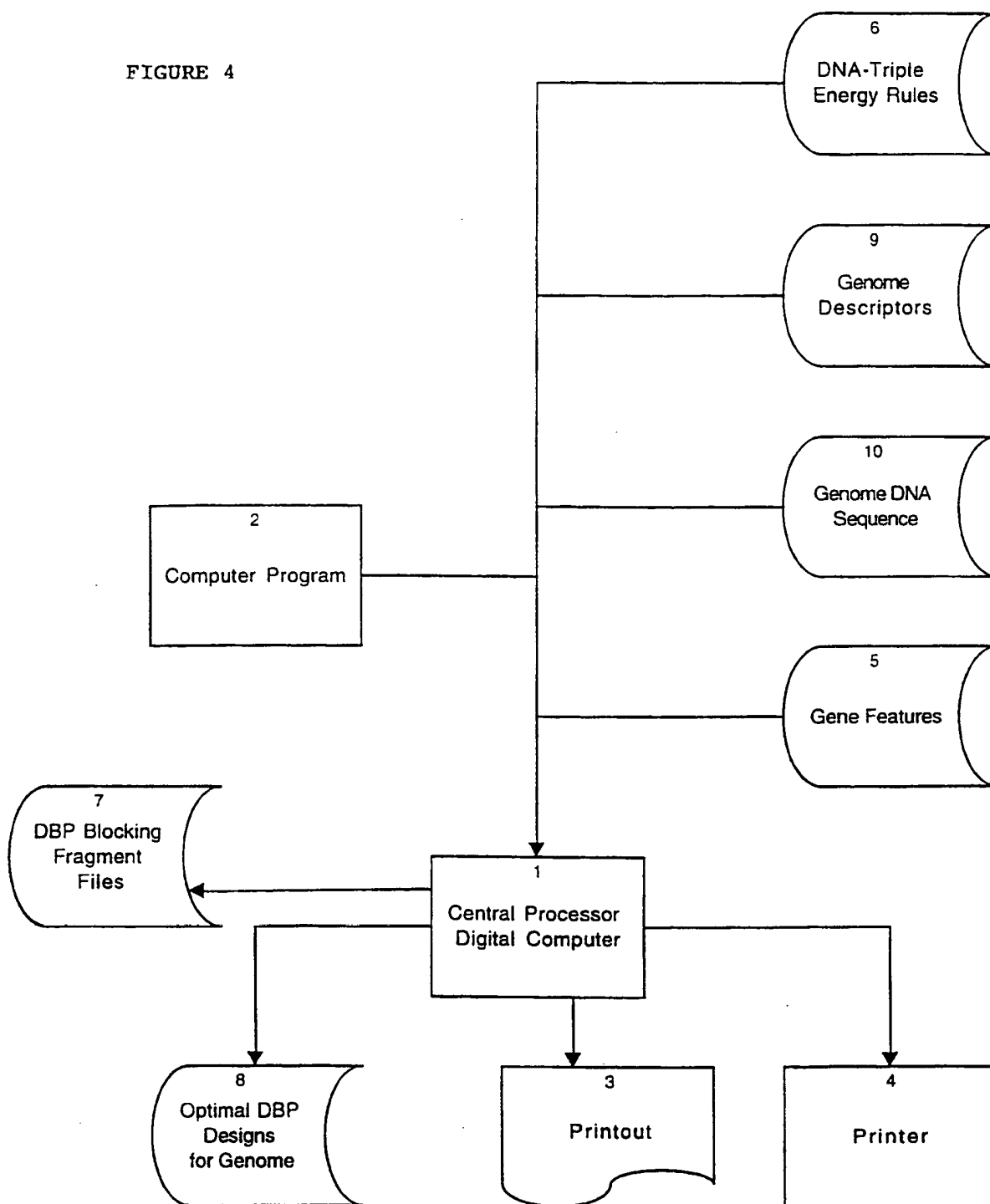


FIGURE 5

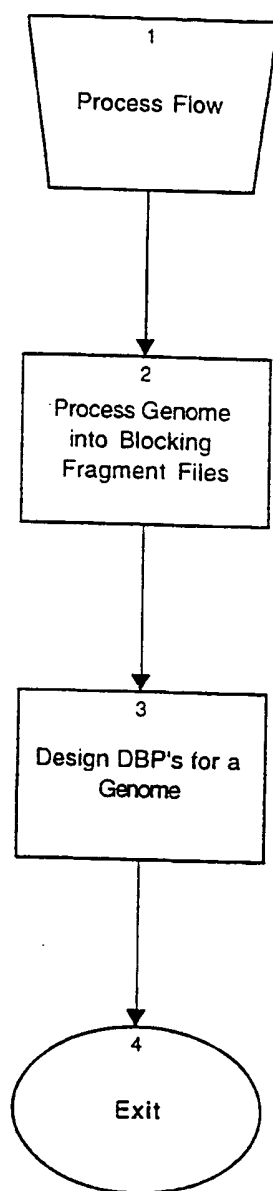
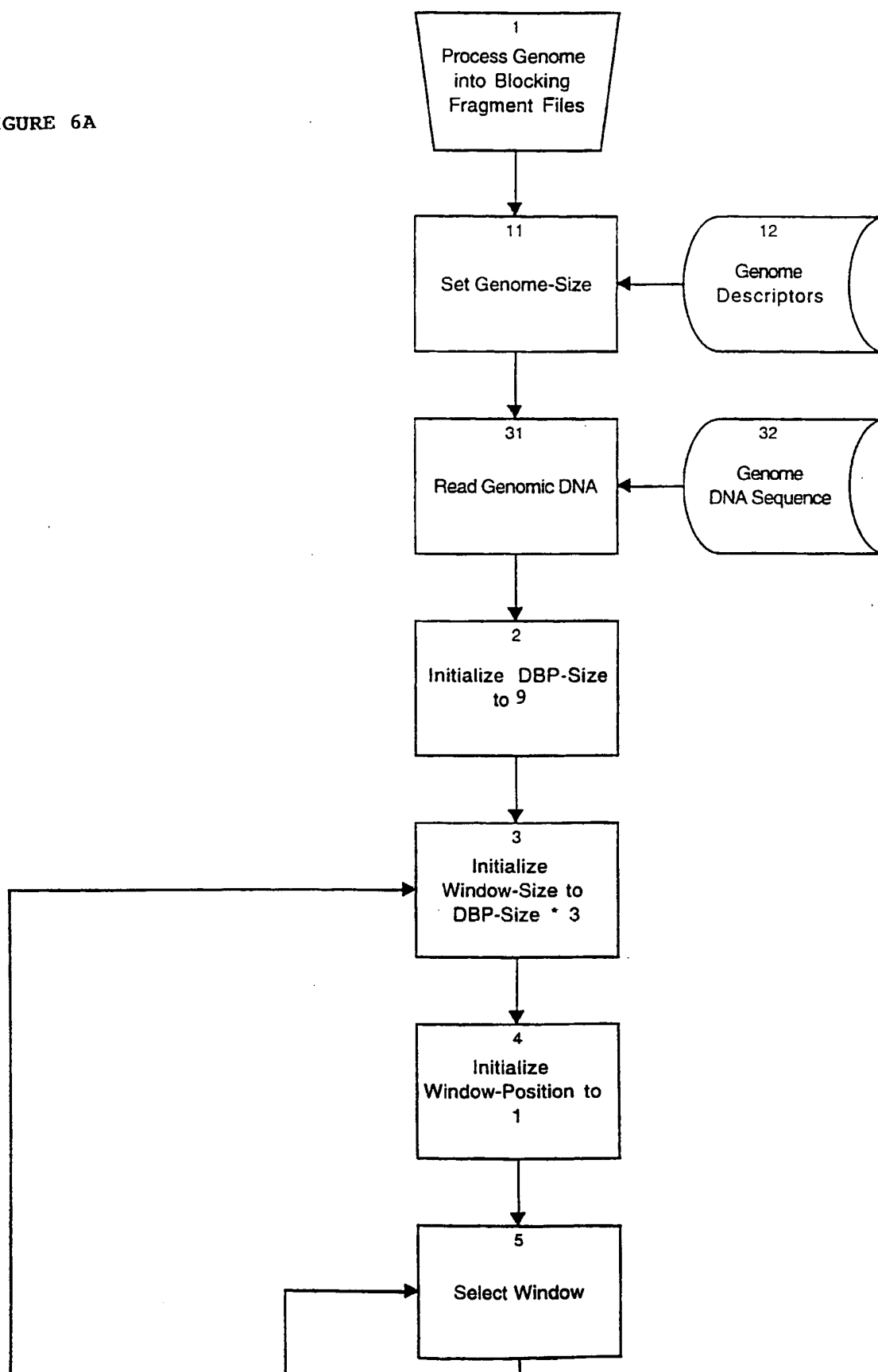


FIGURE 6A



7/26

FIGURE 6B

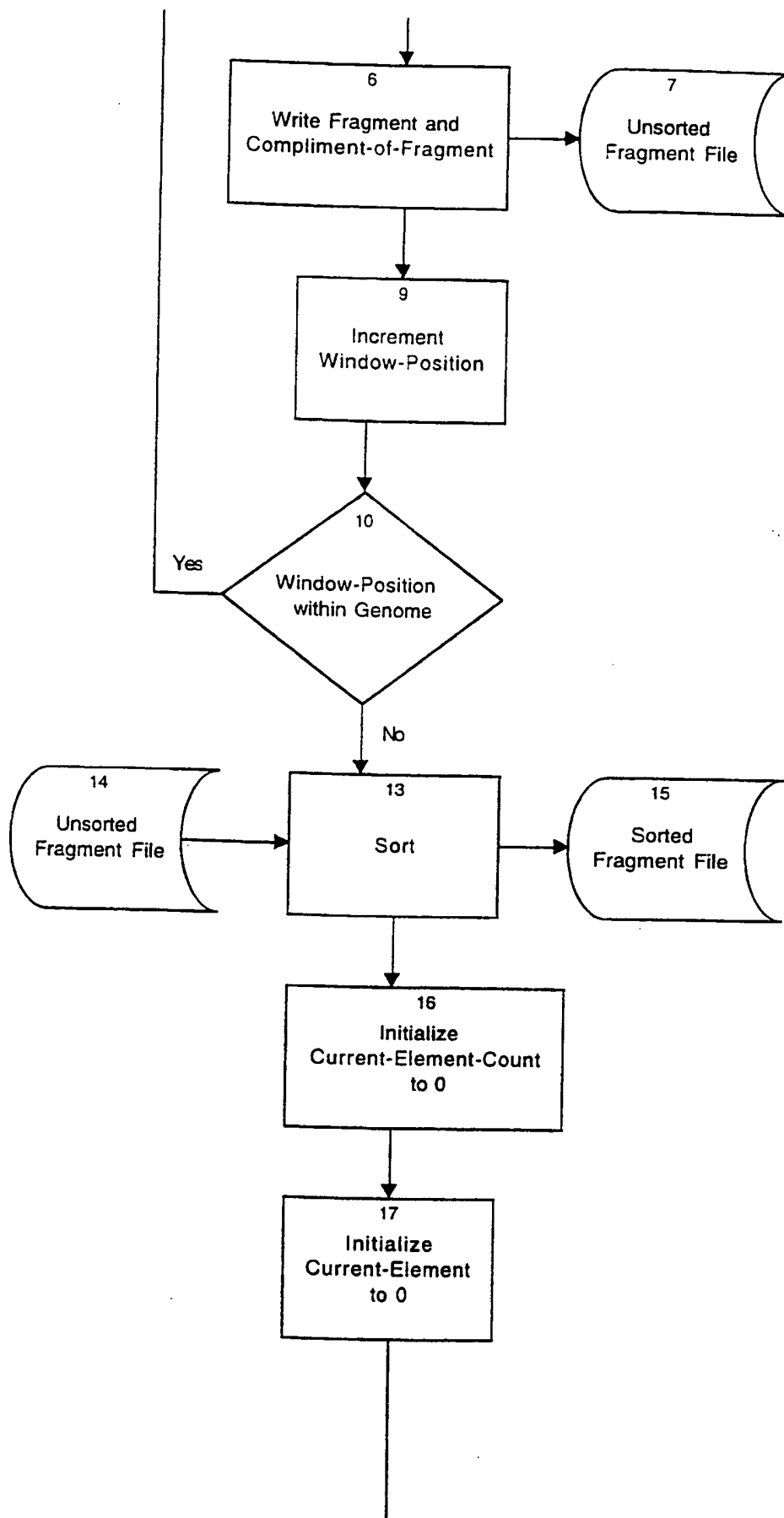
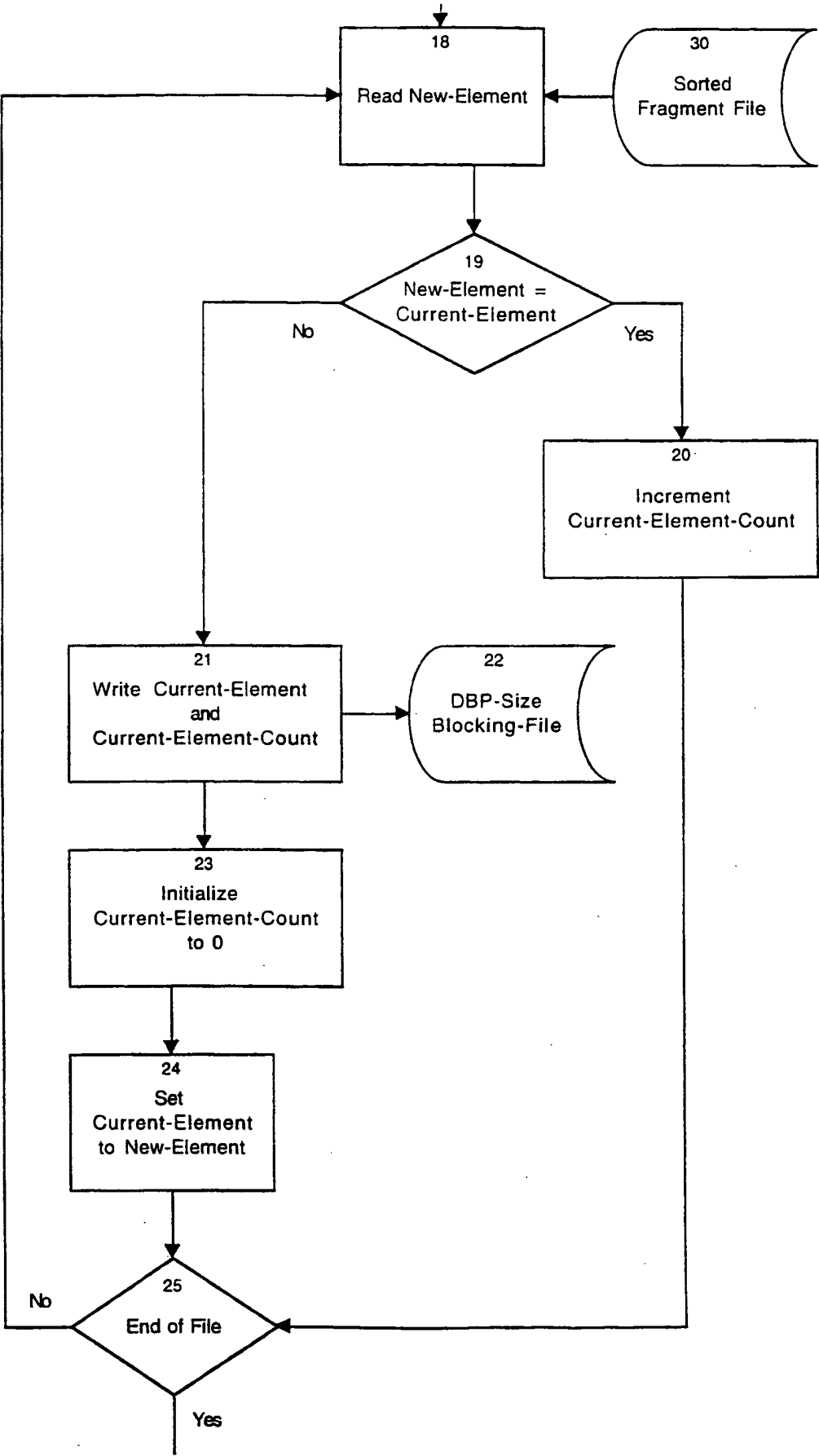


FIGURE 6C



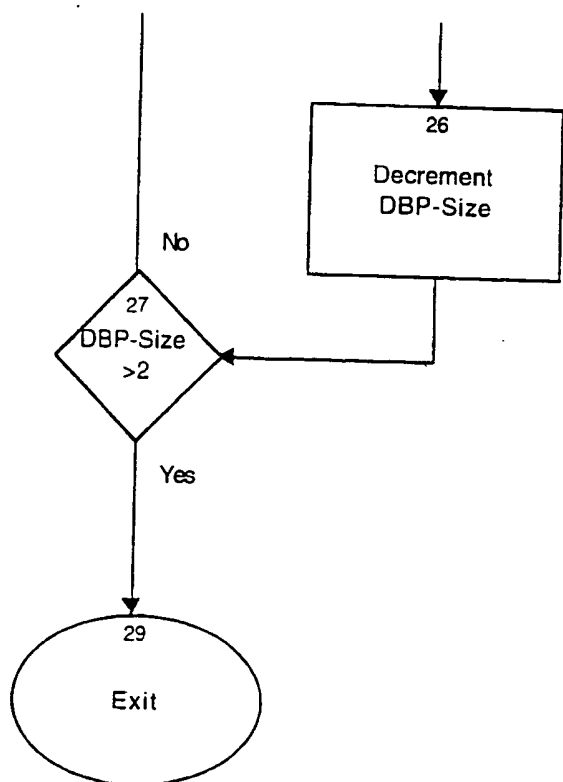
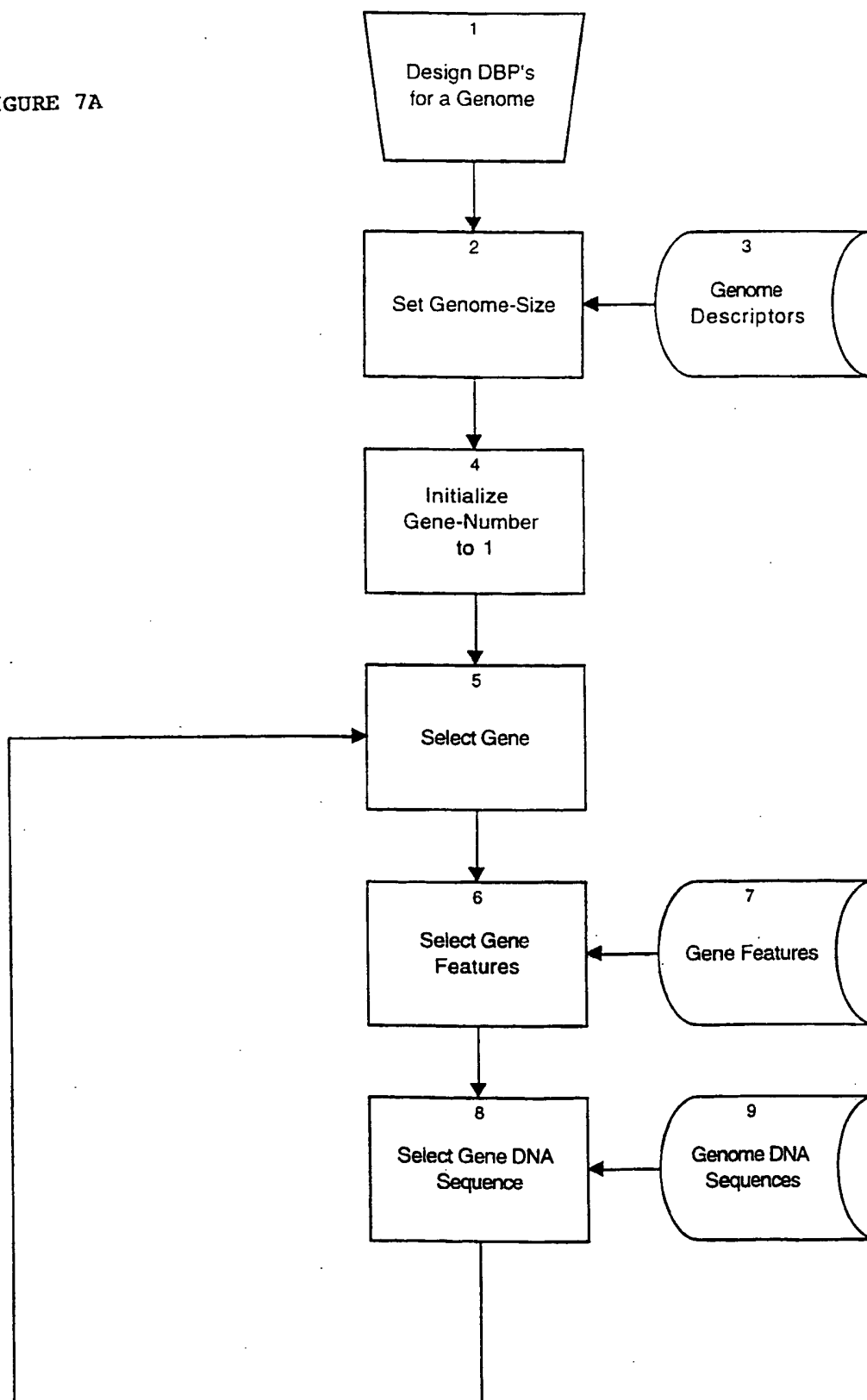


FIGURE 6D

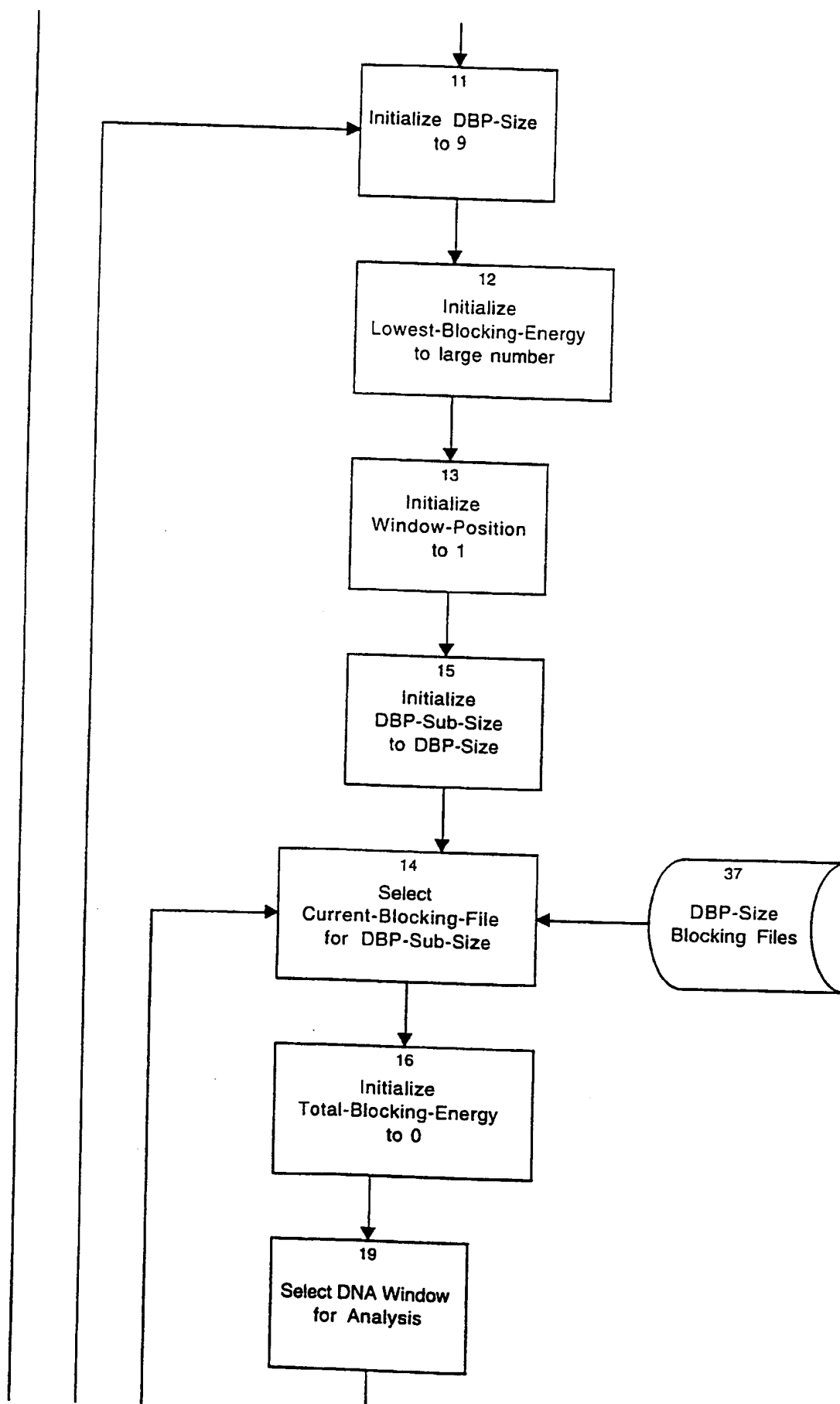
10 / 26

FIGURE 7A



11 / 26

FIGURE 7B





12 / 26

FIGURE 7C

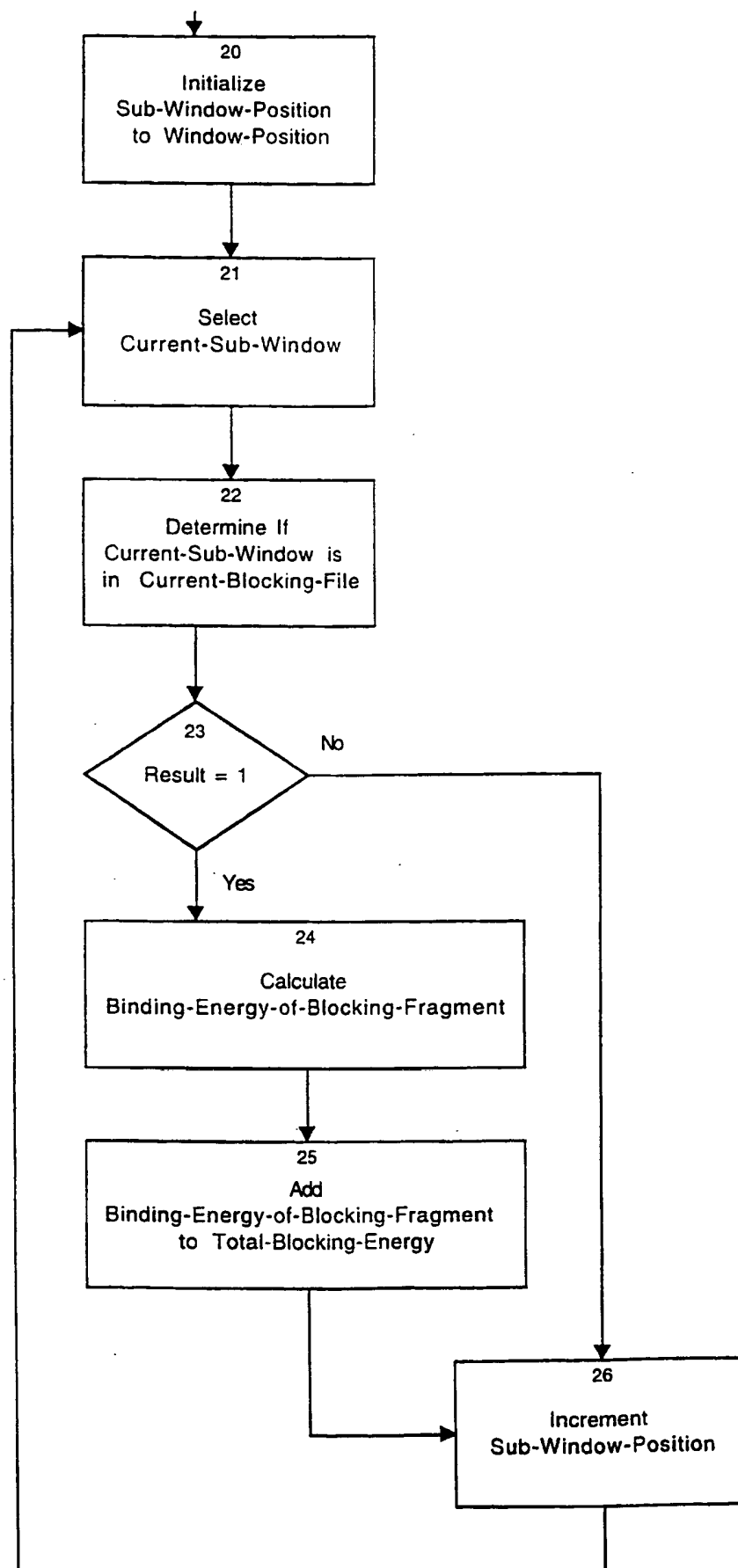
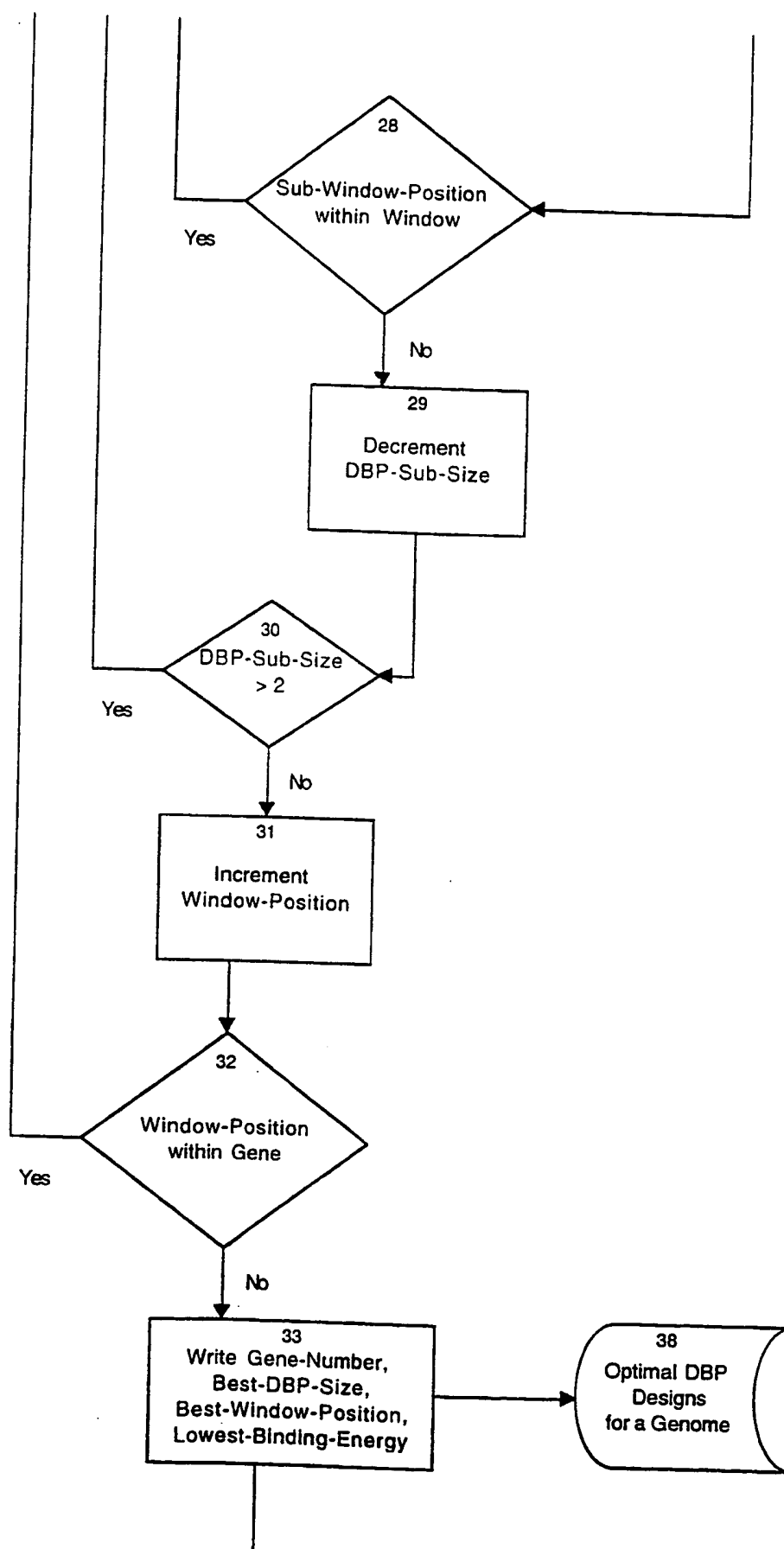


FIGURE 7D



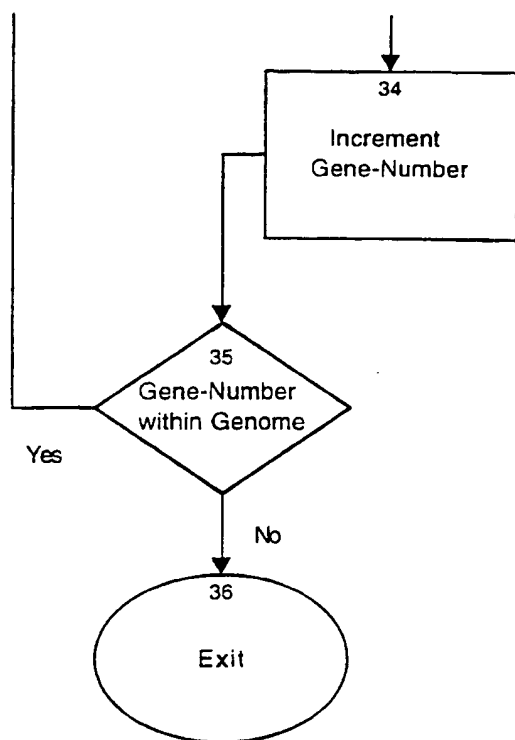


FIGURE 7E

FIGURE 8

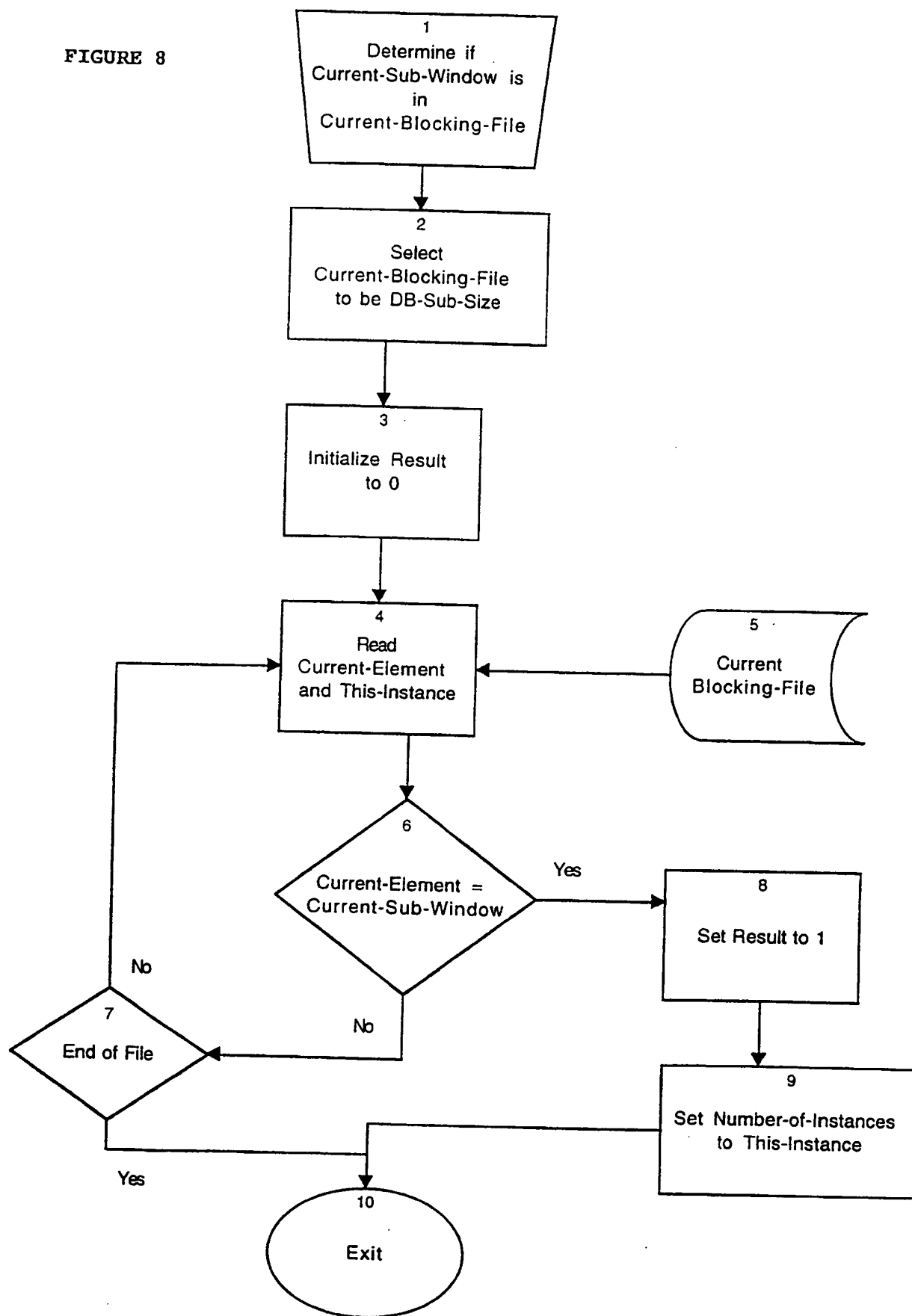
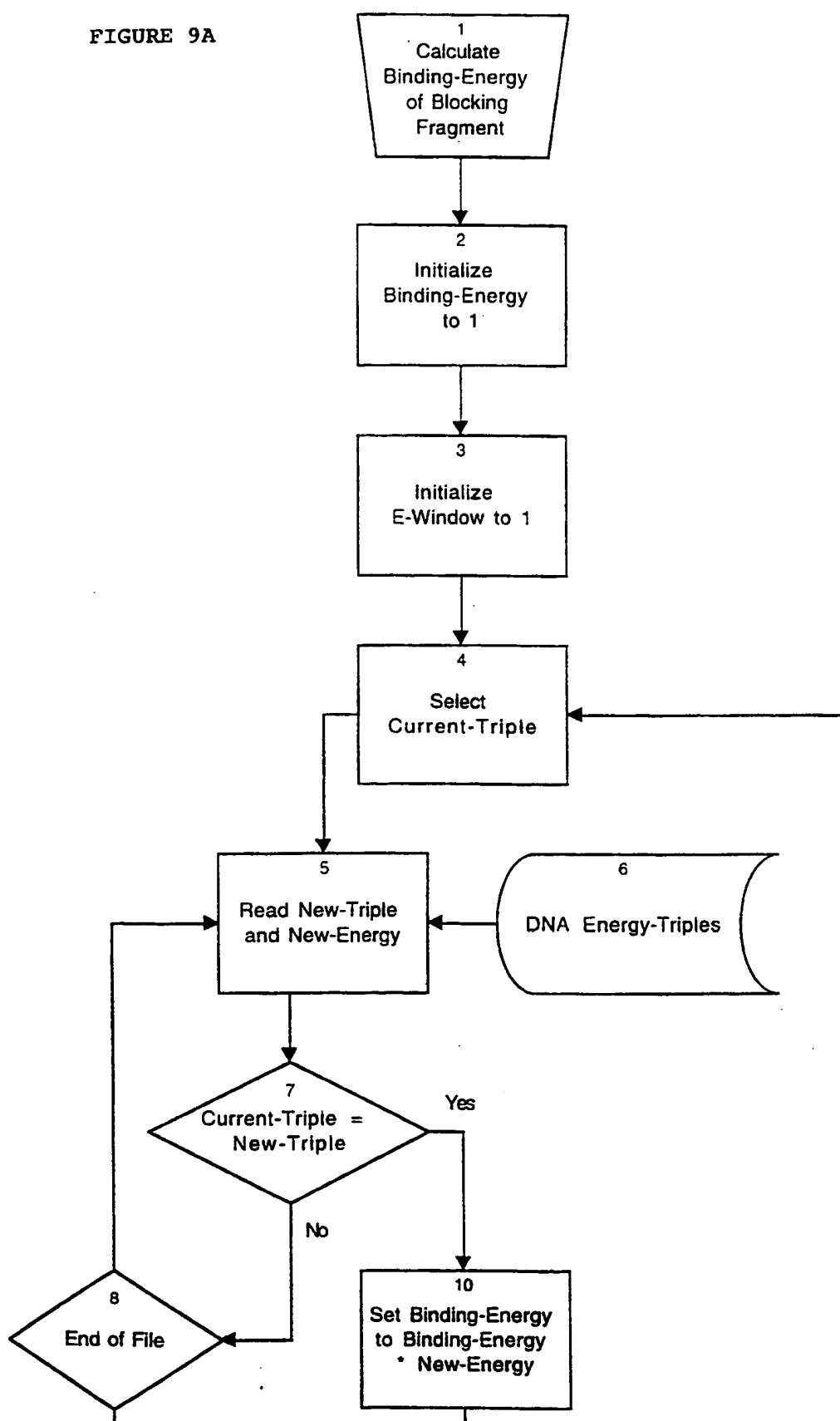


FIGURE 9A



17 / 26

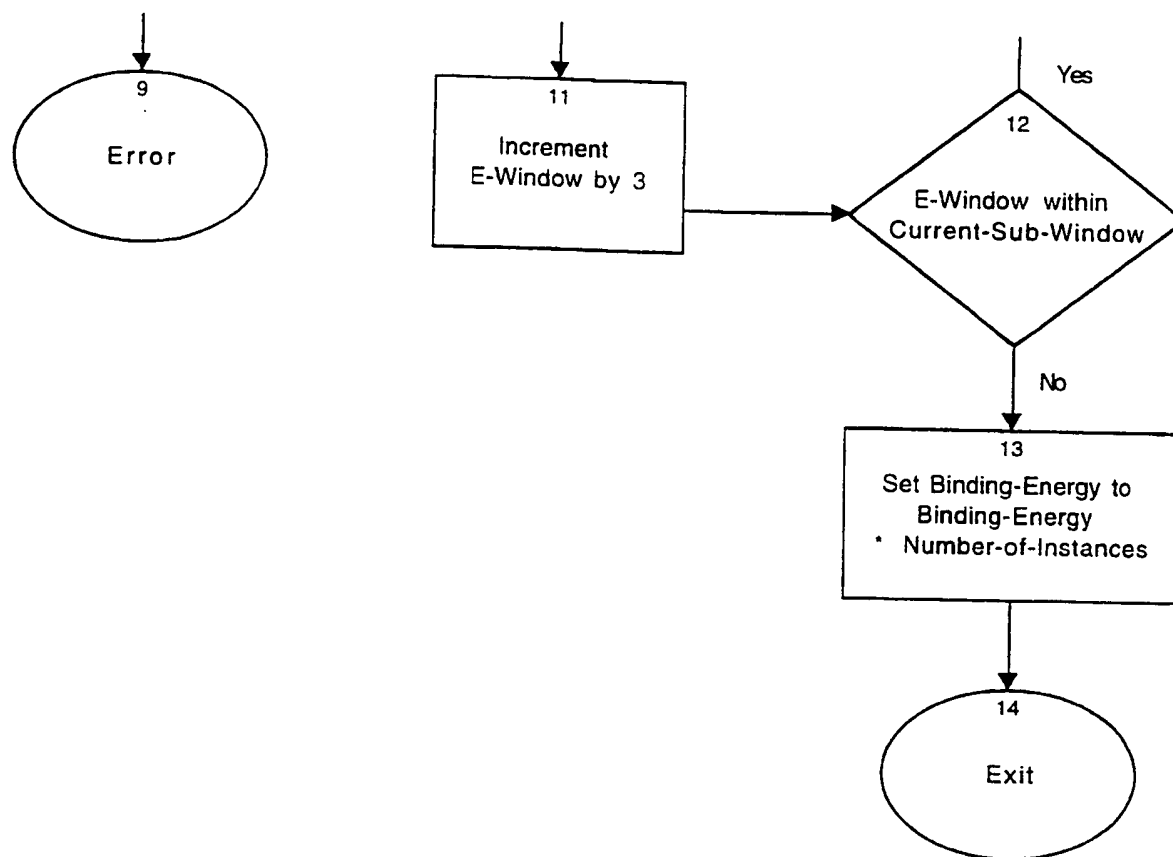


FIGURE 9B

FIGURE 10

Yeast Chromosomes 1 & 2 - Coding Sequence - 1/21/99 -  
Distribution of Strength Order

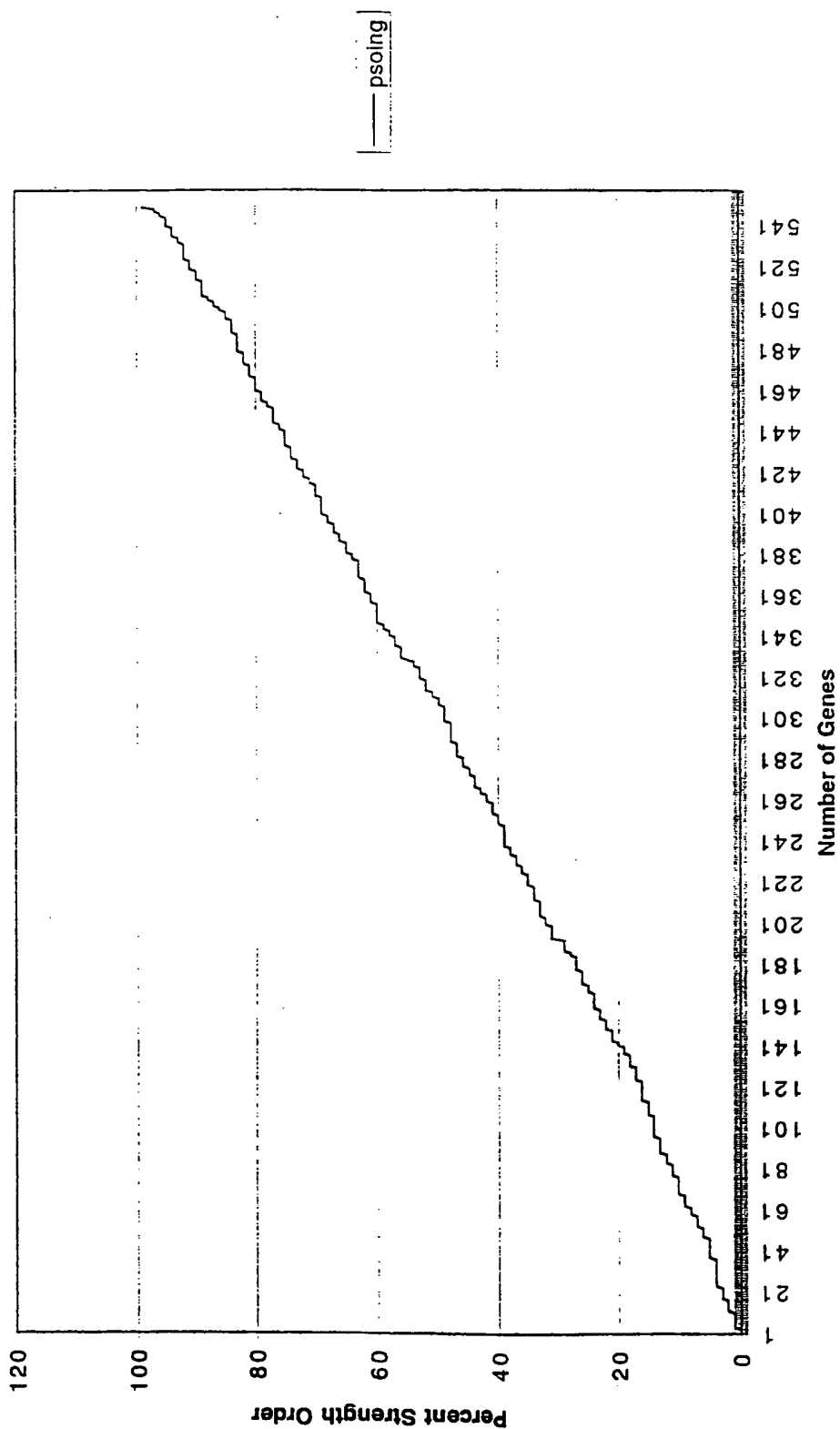


FIGURE 11

# Yeast Chromosomes 1 & 2 - Coding Sequence - 1/21/99 - Distribution of Binding Energy

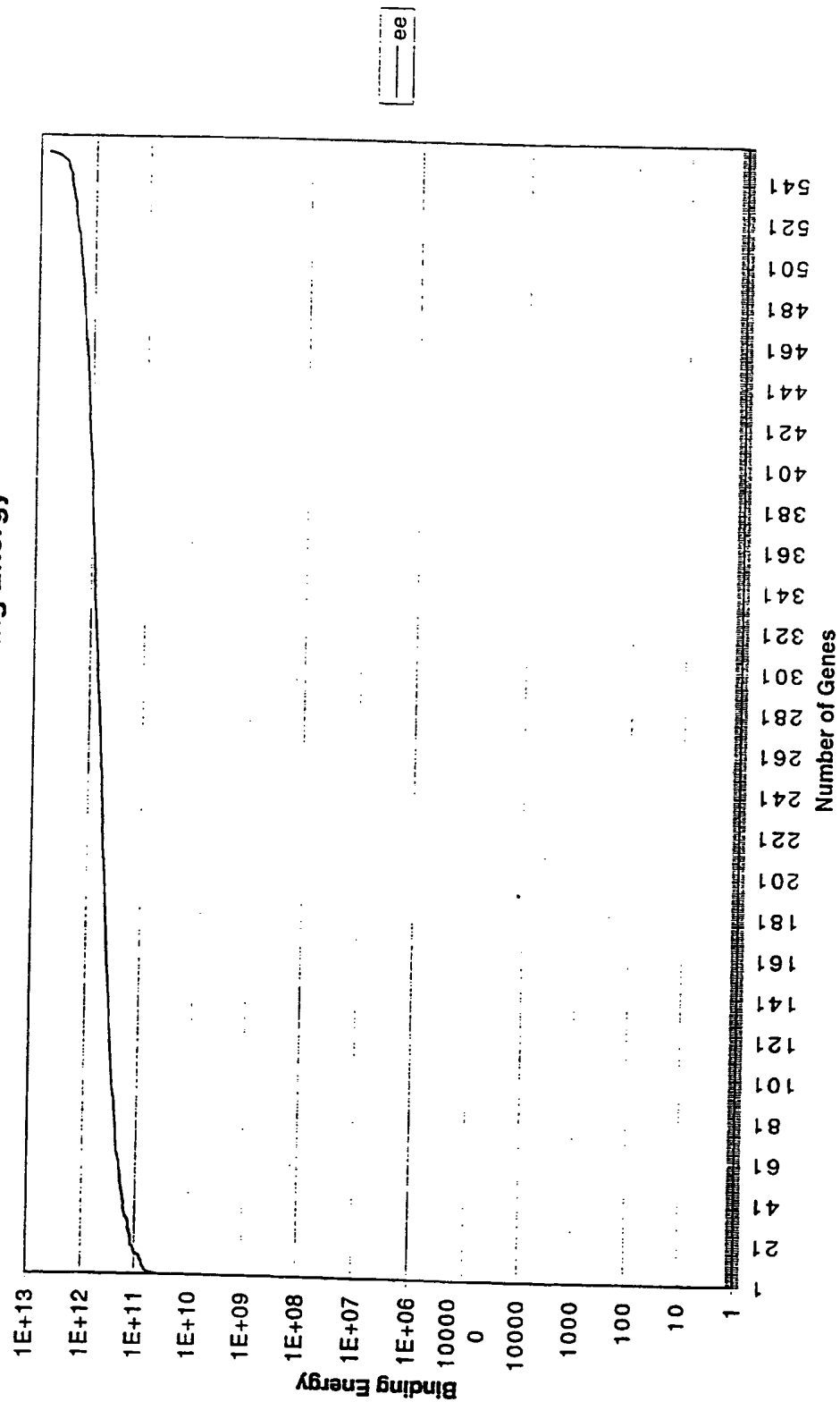




FIGURE 12

**Yeast Chromosomes 1 & 2 - Coding Sequence - 1/21/99 -  
Distribution of Sub-Site Binding Energy**

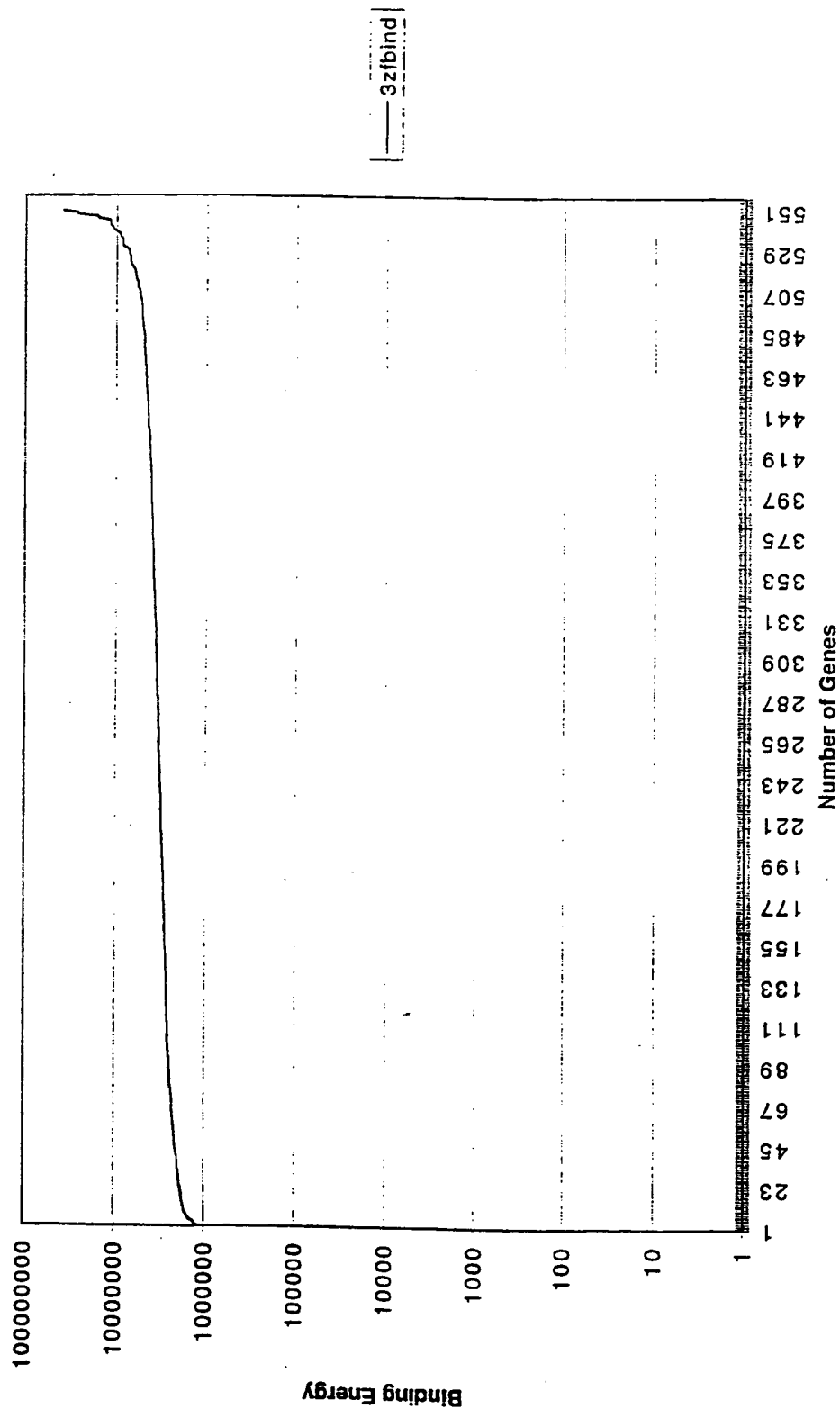


FIGURE 13

**Yeast Chromosomes 1 & 2 - Coding Sequence - 1/21/99 -  
Distribution of Sub-Site Binding Energy**

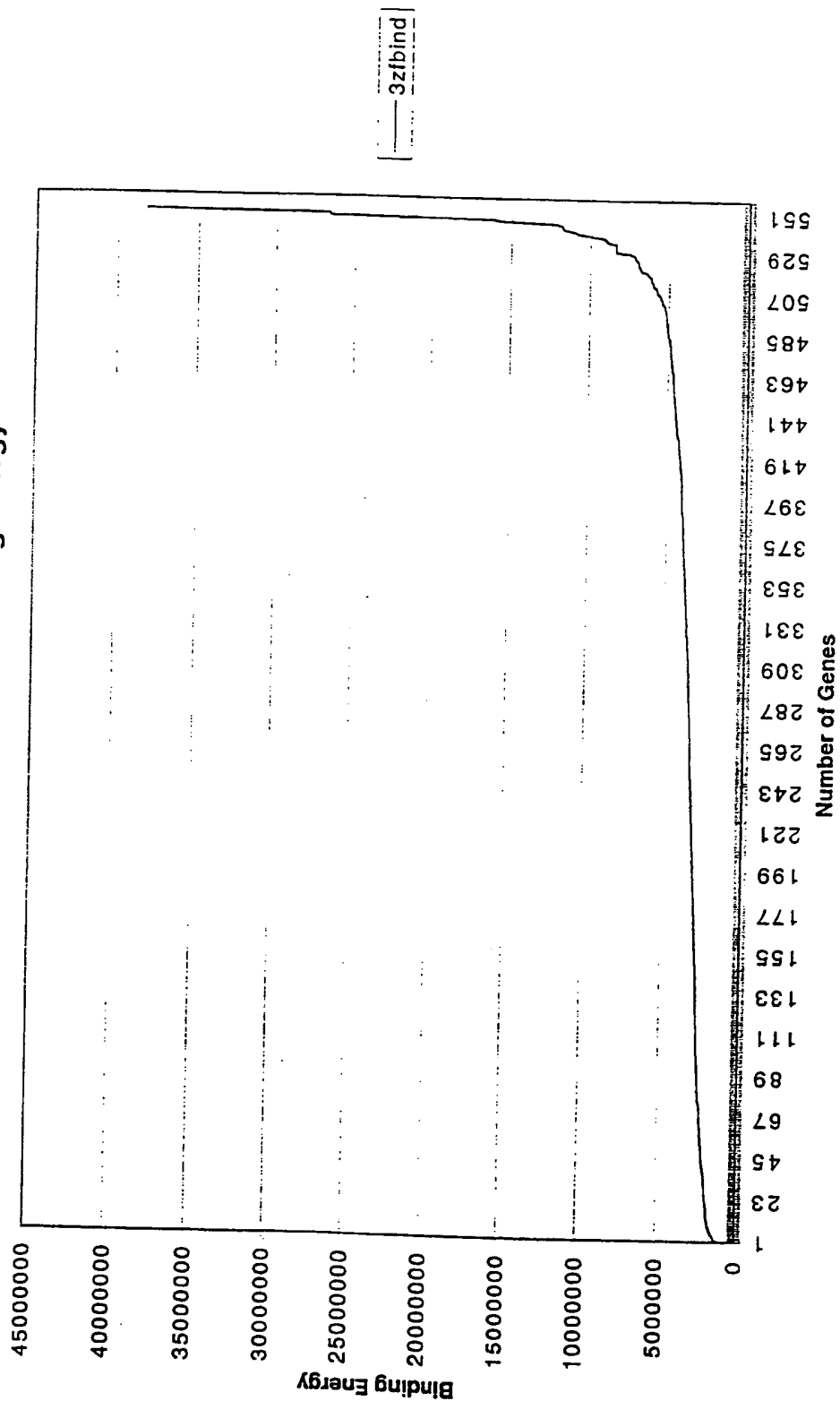


FIGURE 14

Yeast Chromosomes 1 & 2 - Coding Sequence - 1/21/99 -  
Ratio of Binding Energy to Sub-Site Binding Energy

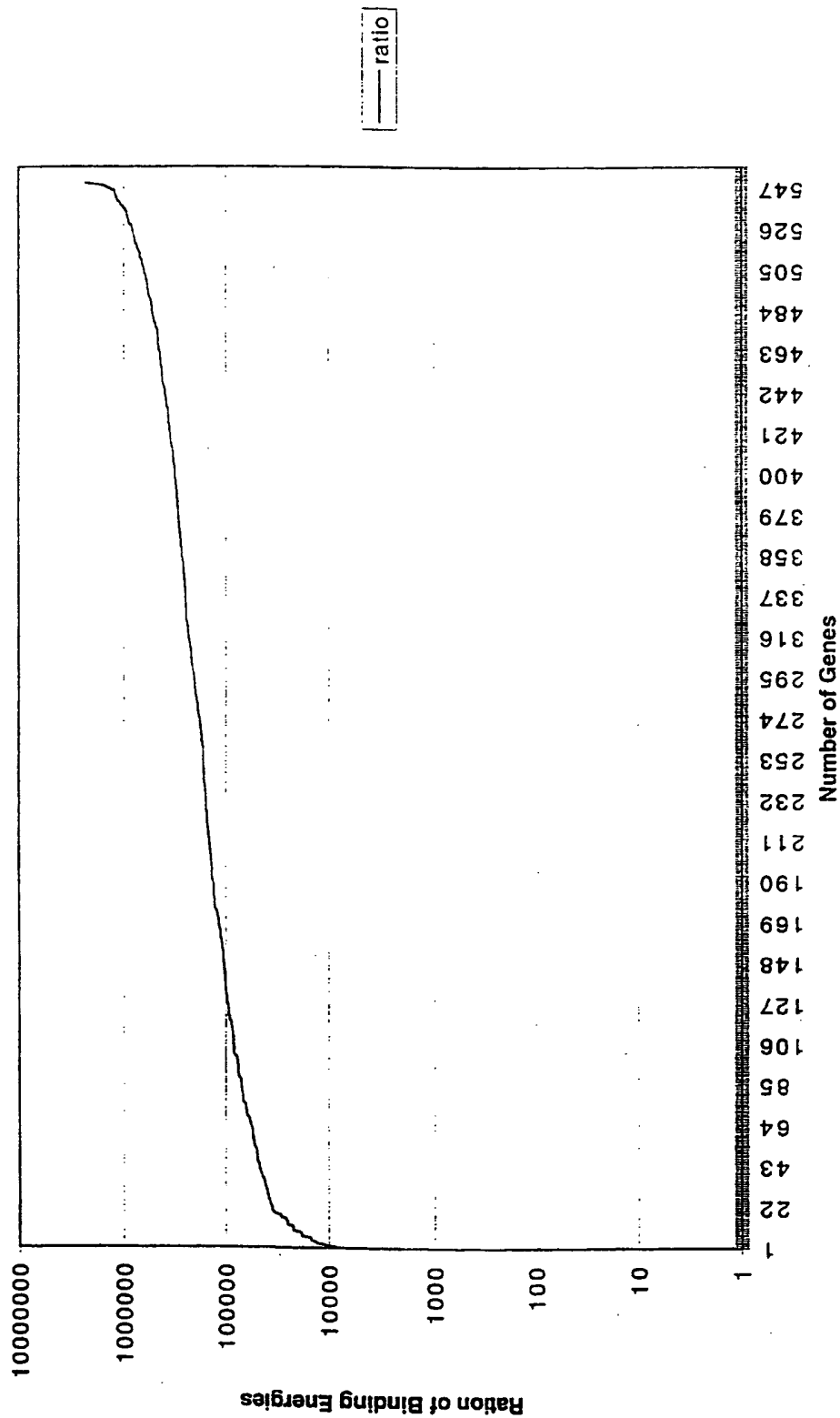


FIGURE 15

Gene YAR073 - Yeast Promoter Sequence - 2/6/99 -  
Spurious Binding Energy

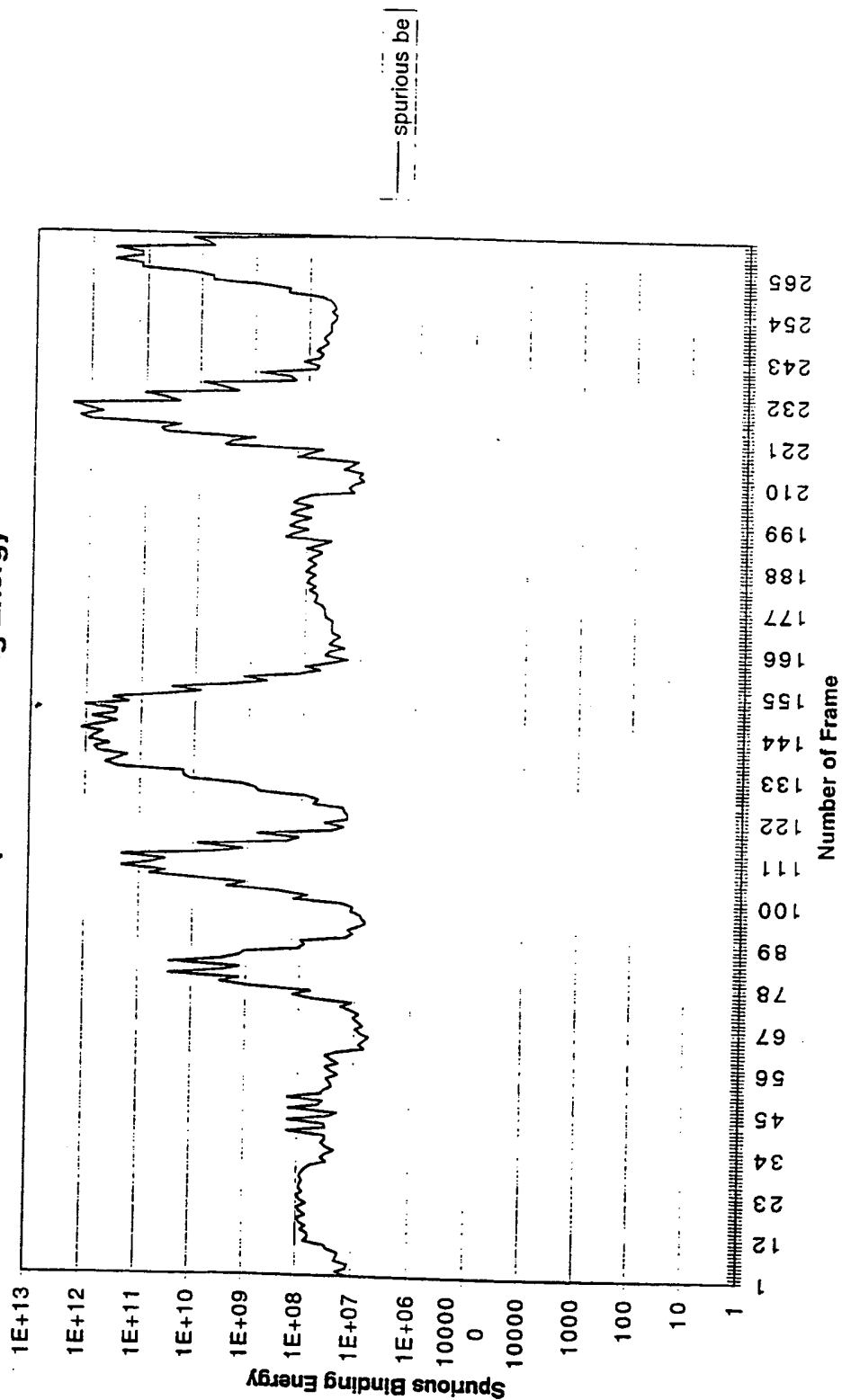


FIGURE 16

Gene YAR073 - Yeast Promoter Sequence - 2/6/99 -  
BE/SBE Ratio

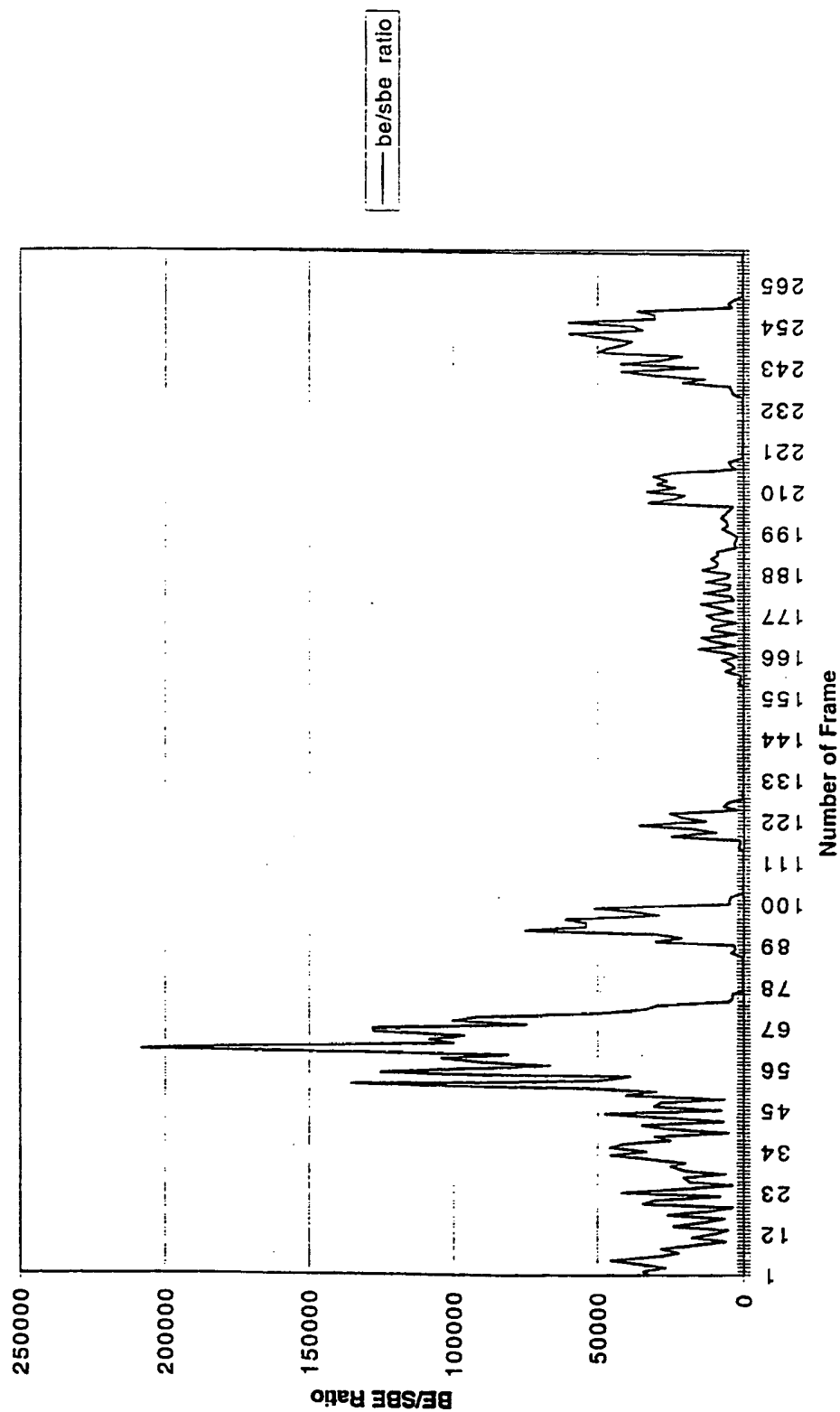


FIGURE 17

Worm Chromosome 1 - Genes 1-100 - Coding Sequence - 2/17/99 -  
Distribution of DBP Sizes

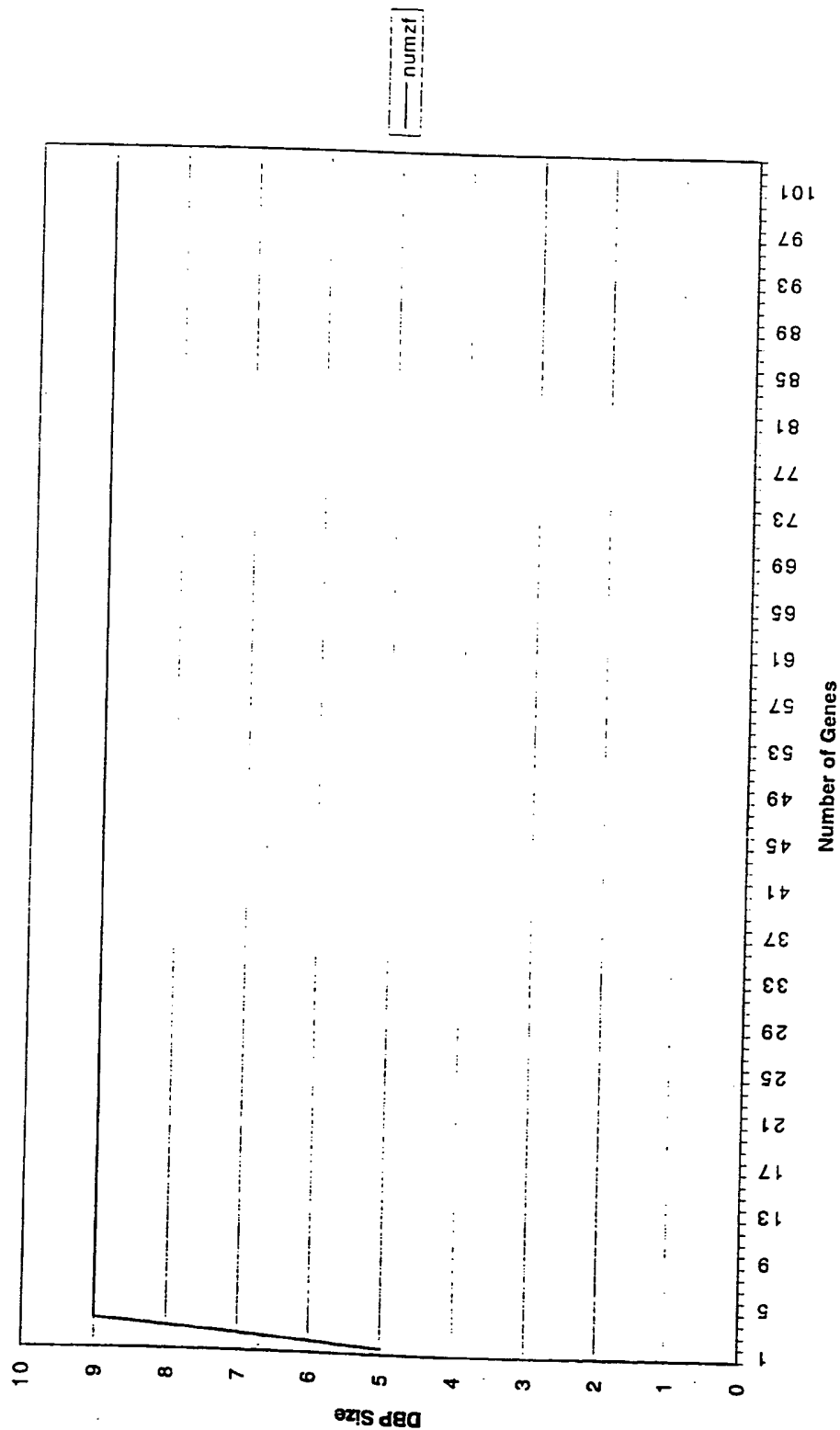


FIGURE 18

Worm Chromosome 1 - Genes 1-100 - Coding Sequence - 2/17/99  
Ratio of Binding Energy to Sub-Site Binding Energy

